

Exhibit B

DESMARAIS LLP

www.desmaraisllp.com

NEW YORK
SAN FRANCISCO
WASHINGTON, DC

Jonas R. McDavit
New York
Direct: 212-351-3425
jmcavit@desmaraisllp.com

May 21, 2021

Via Email

Michael Yang
SVP, General Counsel and Corporate Secretary
Proofpoint, Inc.
925 West Maude Avenue
Sunnyvale, CA 94085
miyang@proofpoint.com

Re: ZapFraud, Inc.

Dear Mr. Yang:

As you are aware, ZapFraud, Inc. (“ZapFraud”) owns a number of active and pending U.S. patents relating to email fraud prevention and detection. ZapFraud was founded by noted electronic security pioneer, Dr. Markus Jakobsson. Dr. Jakobsson is also the current president of ZapFraud.

ZapFraud’s patents are fundamental to the detection of email fraud and the protection of email users. Proofpoint, Inc. (“Proofpoint”) would benefit from a patent license from ZapFraud for the reasons outlined below.

ZapFraud has been studying its patent portfolio, as well as specific Proofpoint product offerings. ZapFraud believes that Proofpoint uses functionality covered by certain ZapFraud patents, and requires a patent license from ZapFraud. For example, Proofpoint infringes claim 1 of U.S. Patent No. 9,245,115 (the “’115 Patent”, see attachment (1)) by providing email security products and services, such as Proofpoint Email Protection, that classify received electronic communications based at least in part on evaluating the content of the aforementioned electronic communications. Proofpoint advertises, markets, and sells its Proofpoint Email Protection product throughout the United States, e.g., <https://www.proofpoint.com/us/products/email-security-and-protection/email-protection>. Please also see the Proofpoint Email Protection product data sheet (attachment (2)) and Proofpoint-authored blog articles attached to this letter (attachments (3) & (4)), which explain the operation of the Proofpoint Email Protection product.

Please advise me when you are available to discuss Proofpoint’s use of the ’115 Patent. You may contact me at your earliest convenience so that we may schedule an appropriate meeting.

DESMARAIS LLP

Michael Yang
May 21, 2021
Page 2

Sincerely,

/s/ Jonas R. McDavit

Jonas R. McDavit

Attachments:

- (1) U.S. Patent No. 9,245,115
- (2) Proofpoint Email Protection Data Sheet
- (3) Understanding BEC Scams: Supplier Invoicing Fraud (Rapp, 2021)
- (4) Using AI to Stop Threats and Reduce Compliance Risk (Pattera, 2020)

cc: Katherine A. Vidal, Winston & Strawn LLP, kvidal@winston.com
David M. Fry, Shaw Keller LLP, dfry@shawkeller.com

Attachment (1)



US009245115B1

(12) **United States Patent**
Jakobsson

(10) **Patent No.:** **US 9,245,115 B1**
(45) **Date of Patent:** **Jan. 26, 2016**

(54) **DETERMINING RISK EXPOSURE AND
AVOIDING FRAUD USING A COLLECTION
OF TERMS**

(71) Applicant: **ZapFraud, Inc.**, Portola Valley, CA
(US)

(72) Inventor: **Bjorn Markus Jakobsson**, Mountain
View, CA (US)

(73) Assignee: **ZapFraud, Inc.**, Portola Valley, CA
(US)

(*) Notice: Subject to any disclaimer, the term of this
patent is extended or adjusted under 35
U.S.C. 154(b) by 0 days.

7,299,261 B1 *	11/2007	Oliver et al.	709/206
7,644,274 B1	1/2010	Jakobsson et al.	
7,899,213 B2 *	3/2011	Otsuka et al.	382/104
7,899,866 B1	3/2011	Buckingham et al.	
7,917,655 B1 *	3/2011	Coomer et al.	709/248
8,010,614 B1 *	8/2011	Musat et al.	709/206
8,131,655 B1	3/2012	Cosoi et al.	
8,255,572 B1 *	8/2012	Coomer	709/248
8,667,069 B1	3/2014	Connelly et al.	
2002/0138271 A1 *	9/2002	Shaw	704/270
2003/0229672 A1 *	12/2003	Kohn	709/207
2003/0236845 A1	12/2003	Pitsos	
2004/0176072 A1	9/2004	Gellens	
2005/0060643 A1 *	3/2005	Glass	G06F 17/241 715/205
2005/0076084 A1	4/2005	Loughmiller et al.	

(Continued)

(21) Appl. No.: **13/765,635**

(22) Filed: **Feb. 12, 2013**

Related U.S. Application Data

(60) Provisional application No. 61/597,972, filed on Feb.
13, 2012, provisional application No. 61/729,991,
filed on Nov. 26, 2012.

(51) **Int. Cl.**
H04L 29/06 (2006.01)
G06F 21/55 (2013.01)

(52) **U.S. Cl.**
CPC **G06F 21/55** (2013.01); **G06F 21/554**
(2013.01); **H04L 63/1408** (2013.01)

(58) **Field of Classification Search**
CPC . H04L 63/126; H04L 12/585; H04L 63/1408;
G06F 21/554
USPC 726/22
See application file for complete search history.

(56) **References Cited**

U.S. PATENT DOCUMENTS

6,161,130 A	12/2000	Horvitz et al.
6,574,658 B1	6/2003	Gabber et al.
7,293,063 B1	11/2007	Sobel

OTHER PUBLICATIONS

"Ahonen-Myka, Helena et al., "Finding Co-Occurring Text Phrases by
Combining Sequence and Frequent Set Discovery," Proceedings of
16th International Joint Conference on Artificial Intelligence IJCAI-
99 Workshop on Text Mining: Foundations, Techniques, and Appli-
cations, (Jul. 31, 1999), 1-9."*

Primary Examiner — Izunna Okeke

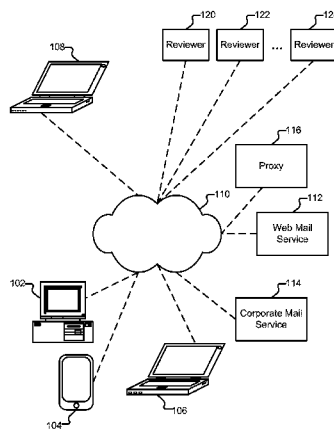
Assistant Examiner — Bryan Wright

(74) *Attorney, Agent, or Firm* — Van Pelt, Yi & James LLP

(57) **ABSTRACT**

Classification of electronic communications includes receiv-
ing an electronic communication, evaluating the received
communication against a collection of terms, and classifying
the received communication based at least in part on the
evaluation. The collection of terms is representative of a
particular strategy of an attacker. The evaluation includes
determining a presence of a portion of the collection of terms
in the electronic communication.

21 Claims, 10 Drawing Sheets



US 9,245,115 B1

Page 2

(56)

References Cited

U.S. PATENT DOCUMENTS

2005/0182735	A1 *	8/2005	Zager et al.	705/67	2007/0271343	A1	11/2007	George et al.	
2005/0188023	A1	8/2005	Doan et al.		2008/0004049	A1	1/2008	Yigang et al.	
2005/0223076	A1	10/2005	Barrus et al.		2008/0046970	A1 *	2/2008	Oliver	G06F 21/554
2005/0235065	A1 *	10/2005	Le et al.	709/238					726/3
2005/0257261	A1 *	11/2005	Shraim	G06Q 10/107	2008/0050014	A1 *	2/2008	Bradski et al.	382/159
				726/22	2008/0141374	A1 *	6/2008	Sidirolou et al.	726/23
2006/0004772	A1 *	1/2006	Hagan et al.	707/10	2008/0276315	A1 *	11/2008	Shuster	H04L 63/1408
2006/0026242	A1 *	2/2006	Kuhlmann et al.	709/206					726/22
2006/0031306	A1	2/2006	Haverkos		2008/0290154	A1 *	11/2008	Barnhardt et al.	235/379
2006/0053490	A1 *	3/2006	Herz et al.	726/23	2010/0030798	A1	2/2010	Kumar et al.	
2006/0149821	A1	7/2006	Rajan et al.		2010/0115040	A1	5/2010	Sargent et al.	
2006/0168329	A1	7/2006	Tan et al.		2010/0145900	A1 *	6/2010	Zheng et al.	706/52
2006/0195542	A1 *	8/2006	Nandhra	709/207	2010/0313253	A1 *	12/2010	Reiss	G06F 21/51
2006/0224677	A1 *	10/2006	Ishikawa et al.	709/206					726/7
2006/0259558	A1	11/2006	Yen		2011/0087485	A1 *	4/2011	Maude et al.	704/9
2006/0265498	A1	11/2006	Turgeman et al.		2011/0191847	A1	8/2011	Davis et al.	
2007/0101423	A1	5/2007	Oliver et al.		2012/0227104	A1 *	9/2012	Sinha et al.	726/22
2007/0192169	A1 *	8/2007	Herbrich et al.	705/10	2013/0081142	A1	3/2013	McDougal et al.	
2007/0198642	A1	8/2007	Malik		2013/0083129	A1 *	4/2013	Thompson	B41J 11/002
									347/51
					2013/0346528	A1	12/2013	Shinde et al.	
					2014/0250506	A1	9/2014	Hallam-Baker	

* cited by examiner

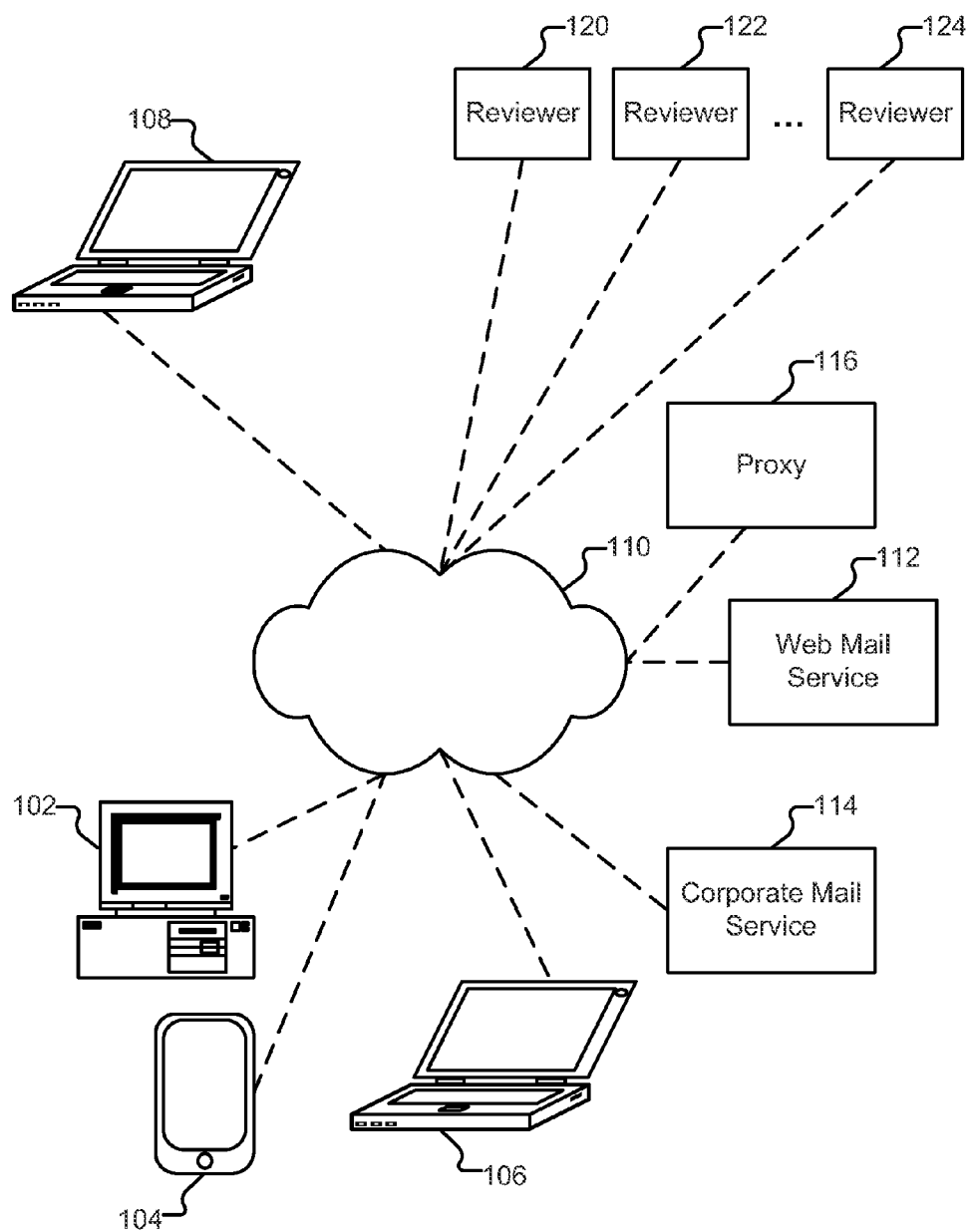


FIG. 1

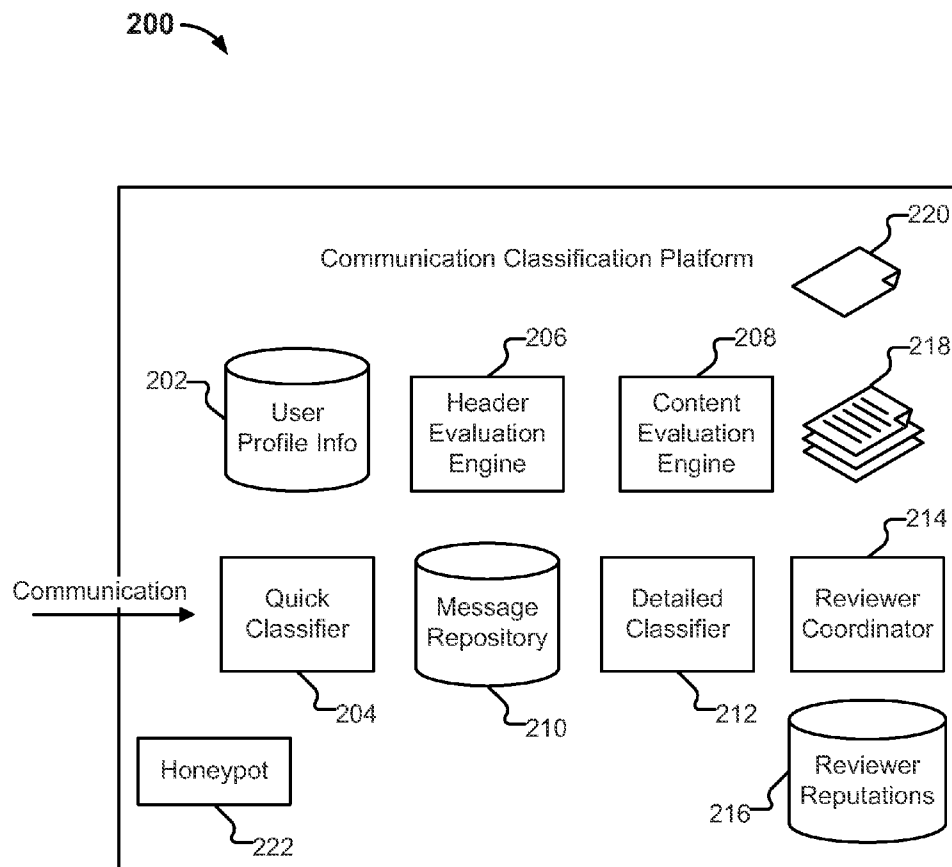


FIG. 2

U.S. Patent

Jan. 26, 2016

Sheet 3 of 10

US 9,245,115 B1

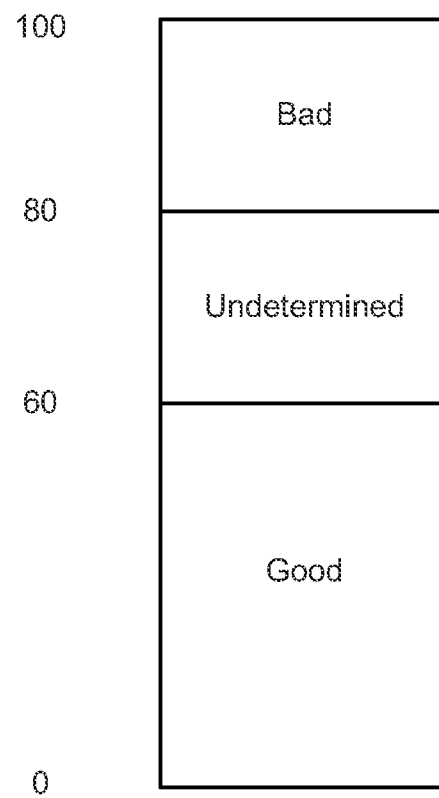


FIG. 3

U.S. Patent

Jan. 26, 2016

Sheet 4 of 10

US 9,245,115 B1

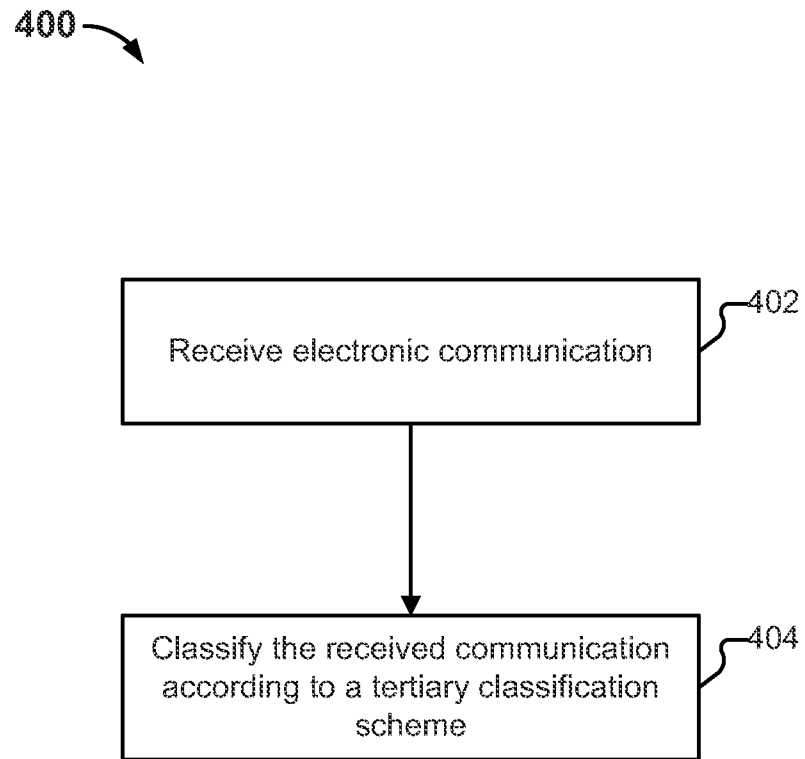
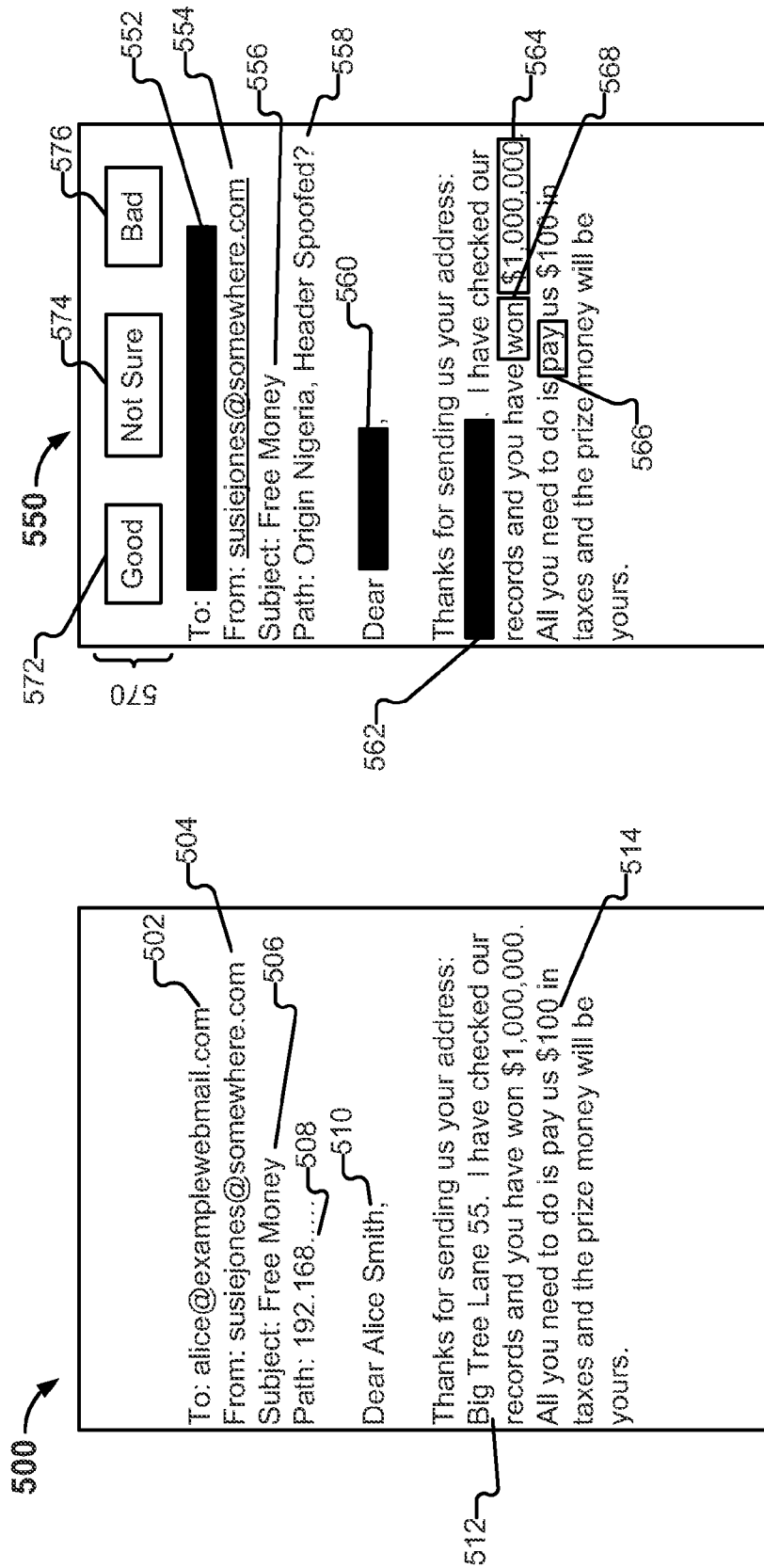


FIG. 4



U.S. Patent

Jan. 26, 2016

Sheet 6 of 10

US 9,245,115 B1

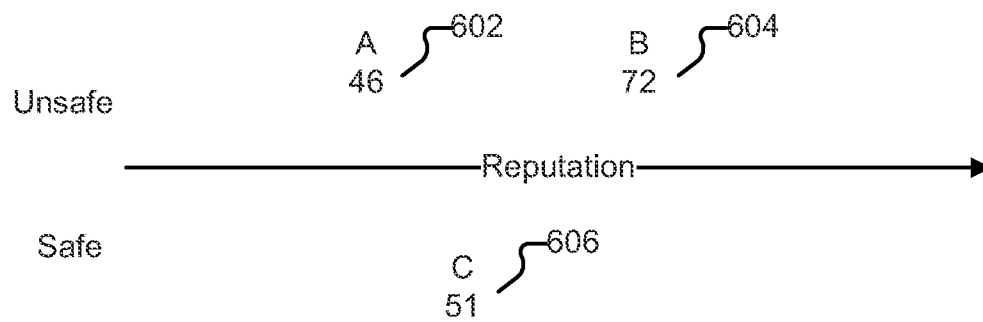
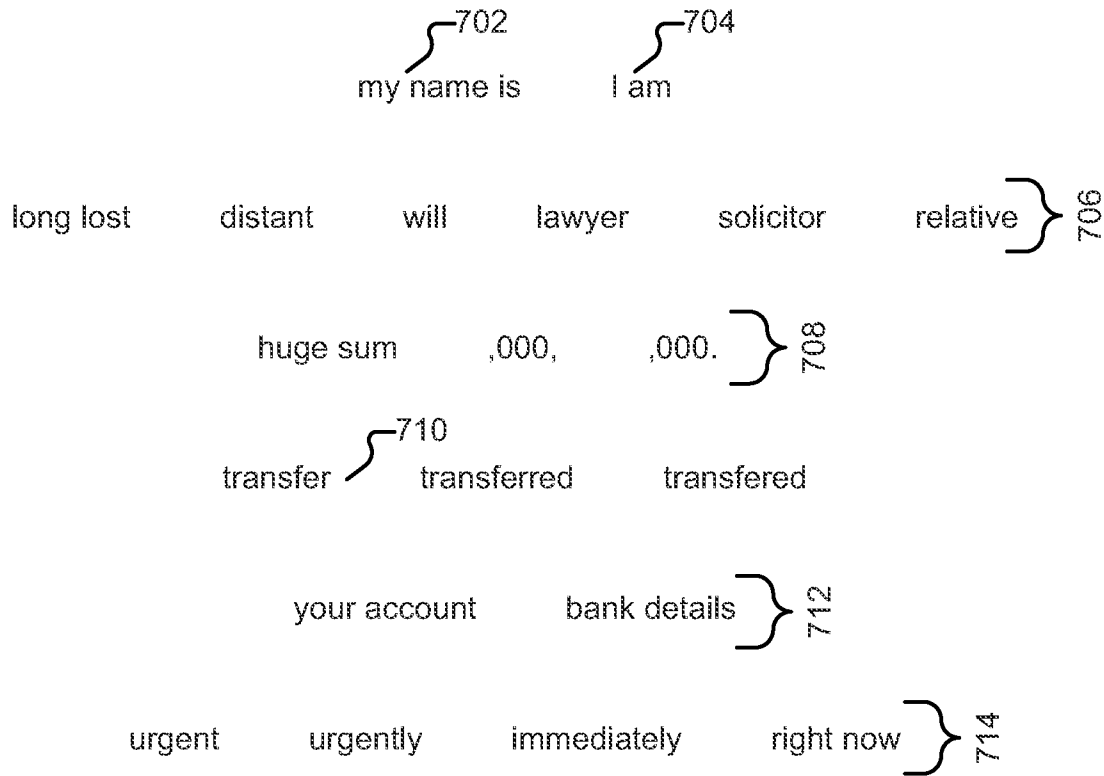


FIG. 6

U.S. Patent**Jan. 26, 2016****Sheet 7 of 10****US 9,245,115 B1**



long lost distant will lawyer solicitor relative } 706
 huge sum ,000, ,000. } 708
 transfer } 710 transferred transfered
 your account bank details } 712
 urgent urgently immediately right now } 714

FIG. 7A

750 →

Hi, my name is Mr. James Smith, and I represent your long lost cousin, who recently died.

He has left you a huge sum of money. I need to transfer the money to your account immediately.

FIG. 7B

U.S. Patent

Jan. 26, 2016

Sheet 8 of 10

US 9,245,115 B1

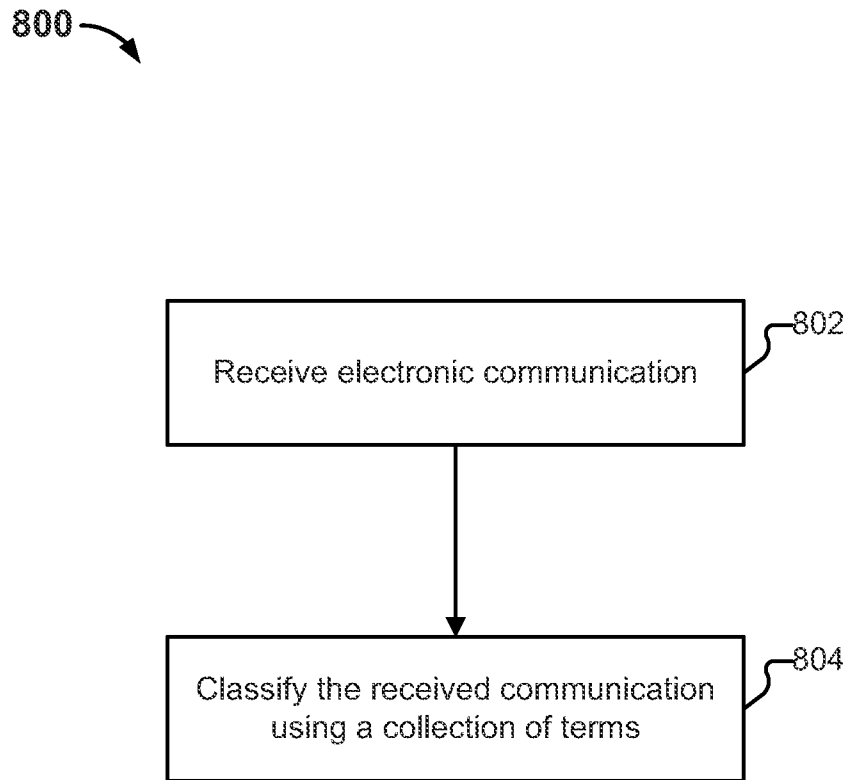


FIG. 8

U.S. Patent

Jan. 26, 2016

Sheet 9 of 10

US 9,245,115 B1

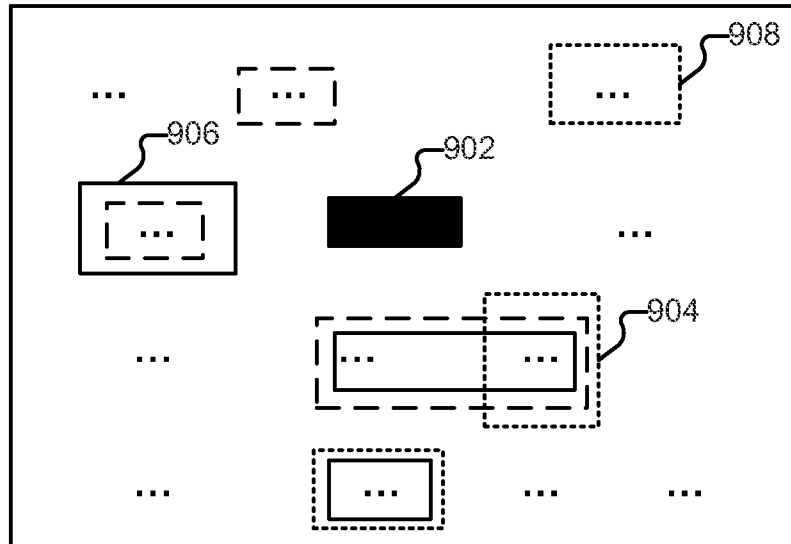


FIG. 9

U.S. Patent

Jan. 26, 2016

Sheet 10 of 10

US 9,245,115 B1

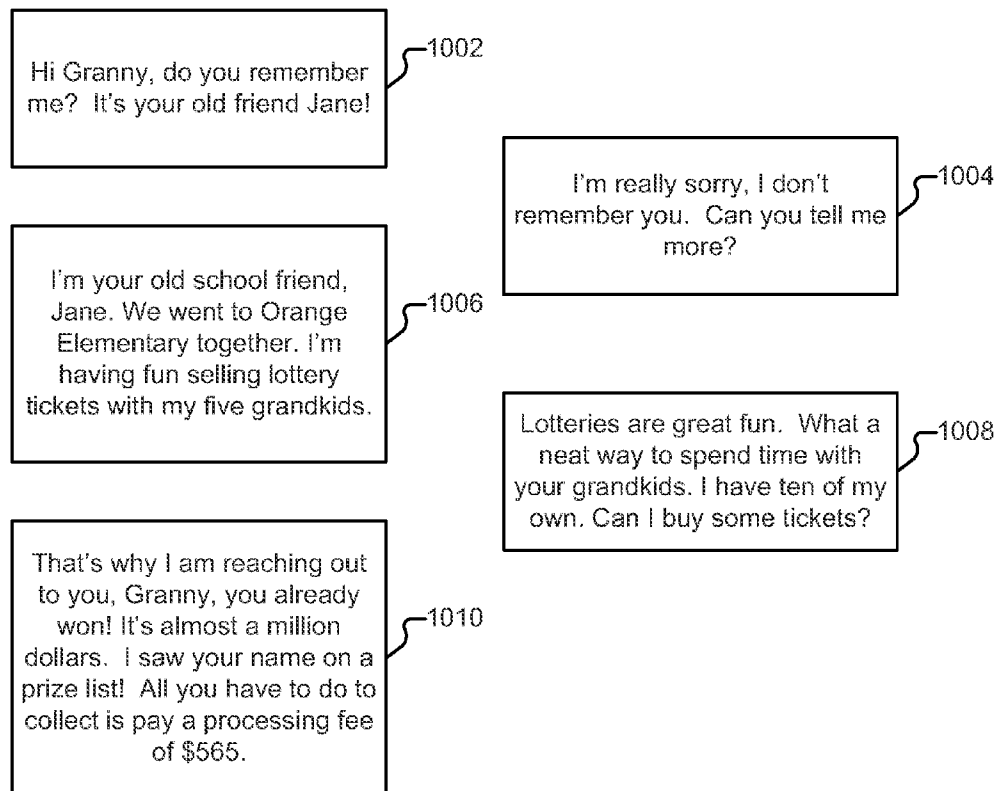


FIG. 10

US 9,245,115 B1

1

DETERMINING RISK EXPOSURE AND AVOIDING FRAUD USING A COLLECTION OF TERMS

CROSS REFERENCE TO OTHER APPLICATIONS

This application claims priority to U.S. Provisional Patent Application No. 61/597,972 entitled PROTECTING COMPUTER AND OTHER DEVICE USERS filed Feb. 13, 2012 which is incorporated herein by reference for all purposes.

This application also claims priority to U.S. Provisional Patent Application No. 61/729,991 entitled DETERMINING RISK EXPOSURE AND AVOIDING FRAUD filed Nov. 26, 2012 which is incorporated herein by reference for all purposes.

BACKGROUND OF THE INVENTION

Individuals increasingly use electronic mail to communicate with one another for personal and business reasons. Unfortunately, unscrupulous individuals can use electronic mail for nefarious purposes, such as to send unwanted advertising email (e.g., SPAM) and perpetrate fraud against victims. Existing techniques for protecting legitimate email users generally rely on the blacklisting of certain terms (e.g., “porn”), where the presence of a blacklisted term in a message automatically results in the classification of the message as SPAM. Such techniques can be readily defeated by the unscrupulous individual. As one example, the unscrupulous individual may use terms that a human would recognize, such as “p0rn” or “p.o.r.n,” but might not appear on a blacklist. More sophisticated approaches are also being undertaken by unscrupulous individuals to defeat blacklisting and other protections. There therefore exists an ongoing need to protect against the victimization of legitimate email users.

BRIEF DESCRIPTION OF THE DRAWINGS

Various embodiments of the invention are disclosed in the following detailed description and the accompanying drawings.

FIG. 1 illustrates an embodiment of an environment in which users of computer and other device are protected from communications sent by unscrupulous entities.

FIG. 2 depicts an embodiment of a communication classification platform.

FIG. 3 depicts an example of a set of score thresholds used in an embodiment of a tertiary communication classification system.

FIG. 4 illustrates an embodiment of a process for classifying communications.

FIG. 5A illustrates an example of an electronic communication.

FIG. 5B illustrates an example of an interface for classifying an electronic communication.

FIG. 6 depicts an example of a review performed by multiple reviewers.

FIG. 7A depicts an example of a collection of terms.

FIG. 7B shows an example of a fraudulent message that would be detected based on the collection of terms depicted in FIG. 7A.

FIG. 8 illustrates an embodiment of a process for classifying communications.

FIG. 9 illustrates an example of an interface configured to receive feedback usable to create collections of terms.

2

FIG. 10 illustrates an example of such a sequence of messages.

DETAILED DESCRIPTION

The invention can be implemented in numerous ways, including as a process; an apparatus; a system; a composition of matter; a computer program product embodied on a computer readable storage medium; and/or a processor, such as a processor configured to execute instructions stored on and/or provided by a memory coupled to the processor. In this specification, these implementations, or any other form that the invention may take, may be referred to as techniques. In general, the order of the steps of disclosed processes may be altered within the scope of the invention. Unless stated otherwise, a component such as a processor or a memory described as being configured to perform a task may be implemented as a general component that is temporarily configured to perform the task at a given time or a specific component that is manufactured to perform the task. As used herein, the term ‘processor’ refers to one or more devices, circuits, and/or processing cores configured to process data, such as computer program instructions.

A detailed description of one or more embodiments of the invention is provided below along with accompanying figures that illustrate the principles of the invention. The invention is described in connection with such embodiments, but the invention is not limited to any embodiment. The scope of the invention is limited only by the claims and the invention encompasses numerous alternatives, modifications and equivalents. Numerous specific details are set forth in the following description in order to provide a thorough understanding of the invention. These details are provided for the purpose of example and the invention may be practiced according to the claims without some or all of these specific details. For the purpose of clarity, technical material that is known in the technical fields related to the invention has not been described in detail so that the invention is not unnecessarily obscured.

FIG. 1 illustrates an embodiment of an environment in which users of computer and other devices are protected from communications sent by unscrupulous entities. In the environment shown, a user of client device 102 (hereinafter referred to as “Alice”) has an email account provided by web mail service provider 112. Alice visits provider 112’s website with her computer (via one or more networks/network types, depicted collectively in FIG. 1 as a single network cloud 110) to read and write email. Alice also accesses her email account via mobile phone 104. A user of client device 106 (hereinafter referred to as “Bob”) has an email account provided by his employer (i.e., hosted by corporate mail server 114) and also has an account with provider 112.

A nefarious individual (hereinafter referred to as “Charlie”) uses client device 108 to perpetrate fraud on unsuspecting victims. In particular, Charlie is a member of a criminal organization that engages in a variety of email scams. One example scam is a “Nigerian scam” (also referred to herein as a “419 scam”), in which criminals contact unsuspecting Internet users with offers, requests, or opportunities for relationships, with the goal of requesting money from the victims—whether within the initial contact email, but also potentially at a future point (e.g., after multiple communications have been exchanged). Nigerian scams are believed to have originated in Nigeria (hence the name), but are now practiced from many places in the world.

While in principle anybody could fall victim to a 419 scam, some Internet users are more prone to such scams than others,

US 9,245,115 B1

3

and many are repeat victims. A particular individual may be the victim of such a scam due to having a personality or needs that make him/her particularly vulnerable to the scam (e.g., an older person living alone). The individual may also be the victim of the scam due to poor security practices by which these users unwittingly invite abuse by sharing their contact information in a way that allows scammers to find it. Furthermore, the names and contact information of these victims may be sold to other scammers, or reused by a successful scammer, due to the high probability of re-victimization.

A 419 scam tends to rely to a larger extent than other scams on social engineering. In contrast to typical SPAM messages which may contain readily blacklistable terms like “porn,” one of the reasons that a 419 scam message is successful at tricking victims is because it appears to be a legitimate conversational message. Terms frequently present in a 419 scam message, such as “dollars” or “account” are also very prevalent in legitimate email. Further, people who are about to fall for a 419 scam may be unwilling to believe that they are being tricked, because they want to believe in the message that the scammer provides them. This makes it difficult for friends and family to help protect victims, as the victims do not believe they are being victimized. The degree of possible customization of scam messages makes it particularly difficult for existing email filters to provide sufficient protection, as evidenced by the ongoing success of such scams.

Described herein are techniques for protecting vulnerable users from malicious entities such as Charlie. In particular, as will be described in more detail below, communications are examined and classified by a classification platform **200**, which can be operated as a single, standalone device, and can also be at least partially incorporated into a variety of the components shown in FIG. 1.

In addition to protecting against 419 scams, a wide array of other structurally related abuses, such as forms of cyber bullying, abuse by sexual predators, and in general, receipt of inappropriate or threatening information or messages, can be protected against. As will be described in more detail below, depending on the nature of the problems typically facing a given vulnerable user, and the severity of these problems, different configurations can be selected. As one example, platform **200** can determine the extent to which different abuses are covered by different parameter choices for a given user after reviewing some email traffic to and from the user in question, by running for a period of time, and/or by statistical methods that compare the user to similar users using the technology. It is possible to start with one configuration and change to another configuration if the first one is not appropriate, whether, e.g., because it is believed to remove desired communications or because it fails to remove communication that is a risk to the vulnerable user. This can be determined among other things from direct feedback from the protected vulnerable user; by manual or automatic scrutiny of quarantined messages, where it is determined what portion of this traffic was legitimate; and/or by scrutiny of the contents and quantity of the mail that is identified as bad. This provides an opportunity to change the settings over time to make them more appropriate for a given protected user, or to adopt the protection features to a changing problem, as the circumstances and exposure of the protected user change.

The techniques described herein can be incorporated into a variety of systems, in a variety of ways, and in various combinations. For example, Alice’s web browser (e.g., on client **102** or client **104**) can be configured to use a plugin whenever she accesses mail service **112**. The plugin can be configured to perform at least some of the techniques described herein as being performed by platform **200**. As another example, a

4

client-side filter can be included on client device **106** and configured to scan/filter all, or a portion of the incoming/outgoing traffic of the device (e.g., traffic between corporate mail service **114** and/or web mail service **112**, irrespective of what type of mail client is used by Bob to access his mail). In yet other embodiments, a regular mail client (e.g., Microsoft Outlook) is modified to support at least some of the techniques described herein.

The techniques herein can also be provided by service providers (e.g., operating embodiments of platform **200** or configuring their infrastructure to cooperate with embodiments of platform **200**). For example, Alice’s ISP, web mail service **112**, and corporate mail service **114** can each/all provide services in accordance with the techniques described herein. In particular, existing infrastructure provided by the service provider(s) can be adapted to provide at least some of the services described herein, or such services can be provided by one or more separate modules. For example, at least a portion of the functionality of platform **200** can be provided as a gateway (e.g., such that all of the mail of an enterprise is filtered by such a gateway as it arrives/is downloaded). As another example, the functionality of platform **200** can be provided at least partially by a filter (e.g., such that some portion of message processing is performed for free on behalf of an ISP, and any usage above that portion is charged to the ISP based on a contractual agreement). As yet another example, such services can be provided by proxies. The proxies can be controlled by the service providers (e.g., on premises), and can also be provided by a third party as an external service to the service provider. Service providers may opt to provide the services described herein to all users, and can also provide the services on a per-user basis. As one example, Alice could choose to “opt-in” to having web mail service **112** provide her with protections, while Bob could choose to forgo having web mail service **112** provide him with protection. As another example, Alice’s ISP might automatically enroll her in protection services (e.g., based on her age or other demographic information indicative of her being particularly vulnerable, or based on any prior history of Alice having been victimized in an email scam). Bob, having different characteristics, would not automatically be enrolled (but could nonetheless be offered the ability to opt-in).

In some embodiments, a proxy **116** performs scanning/filtering services on behalf of users as a third party subscription service. For example, Alice’s daughter “Eve” can purchase a subscription on behalf of her mother, Alice, or Alice might purchase a subscription for herself. As another example, Bob can purchase a subscription for himself, or his employer might subsidize the subscription on his behalf. The proxy is configured with the user name(s) and password(s) or other credentials for the email accounts to be protected. The usernames/credentials can be collected in a variety of ways. As one example, the user (e.g., Alice) can be explicitly asked to provide them. As another example, the information can be automatically collected on her behalf. The proxy can then access the account(s) periodically, or screen and filter traffic as it is being sent to and from the subscribed user.

In situations such as where services are provided by a third party (e.g., protections for Alice’s account with web mail service **112** are provided by third party proxy **116**), a filter can be used in order to avoid polling the email service provider (e.g., service **112**) too often, which could be considered abusive/disruptive to the operation of service **112**. One way to provide such filter services is by changing the DNS lookup tables associated with the client device, and thereby causing all traffic to be filtered by the proxy as it is being downloaded from web mail service **112** to the client device. Another

US 9,245,115 B1

5

approach is to augment the access functionality on the client device so that proxy **116** is notified when the client device attempts to access the account. As it is being notified, it performs the filtering activity. In yet other embodiments, proxy **116** provides a middleware component to device **102**, where the middleware component catches calls made by the mail reader program (or browser) residing on the client device and then initiates a call to the web mail service **112**. In yet other embodiments, the proxy mimics a client device that is constantly logged in and is polling for updates at a frequent but reasonable rate, such as once every minute.

In various embodiments, combinations of the components described above are used. For example, Alice can be protected both by a plugin or executable installed on client device **102**, and one or more external protection services (e.g., offered by her ISP, by proxy **116**, or by web mail service **112**). In this scenario, Alice will be protected both when she uses her client computer **102** to read her mail, and also when she checks her mail in other places, such as at a library terminal or hotel kiosk. Finally, while the techniques described herein are generally described in conjunction with evaluating email communications, other forms of communications can also be monitored/filtered as applicable. For example, instant messaging clients can monitored (whether at the client, via a proxy, or at a server), and messages being sent to/from the protected user on such services treated in a similar way as is described for the emails, as applicable. SMS/MMS messages are another example of communications that can be screened/managed using the techniques described herein. Other communication technologies can also be monitored and filtered, as applicable. For example, automated voice recognition techniques could be used in conjunction with the screening of voicemail messages (e.g., in conjunction with a service such as Google Voice) or calls, and escalation involving human review could be performed (e.g., with the consent of the callee).

—Communication Classification Platform—

FIG. 2 depicts an embodiment of a communication classification platform. As shown in FIG. 2, platform **200** can comprise a single device, such as standard commercially available server hardware (e.g., with a multi-core processor, 4+ Gigabytes of RAM, and one or more Gigabit network interface adapters) and run a typical server-class operating system (e.g., Linux). Platform **200** can also be implemented using a scalable, elastic architecture and may comprise several distributed components, including components provided by one or more third parties. As explained above, platform **200** (or portions thereof) can be incorporated into a variety of different components depicted in the environment of FIG. 1. As one example, all or portions of platform **200** may be provided by web mail service **112**. As another example, portions of platform **200** may be located on client (or other) devices, such as client device **102** and portions not located on the client device may be omitted, or provided by a third party, as applicable.

In some embodiments, platform **200** includes a database **202** of user profile information. As one example, where proxy **116** implements platform **200**, database **202** could include, for each user of the proxy, the user's username/password information for sites that are proxied. Database **202** can also include information such as the user's credit card information (where the proxy is run as a paid service), contact information, and any user-specific customizations. Examples of such customizations include user-specific whitelists (and any contextual information used to construct those lists, such as temporal information associated with message exchange), scoring thresholds, etc., described in more detail below. As another

6

example, where client device **102** implements platform **200**, database **202** can be implemented as one or more configuration files specific to the user(s) of the device. Further, in some embodiments communications for all users are handled the same way, and database **202** (and/or user-specific profile information) is omitted, or reduced in scope, as applicable.

In some embodiments, when platform **200** receives a communication for processing, the communication is provided to quick classifier **204**. Header evaluation engine **206** evaluates any headers associated with the communication. Examples of information evaluated by engine **206** include: the sender/originator of the communication, the apparent location and IP address of the sender, and the type of domain used by the sender. The header evaluation engine can also evaluate circumstances associated with the communication transmission, such as the time of day it was received, and whether it appears to be a first contact with the user, or a subsequence communication. Content evaluation engine **208** evaluates the content of the communication. As will be described in more detail below, examples of content analysis include analysis based on a "collection of terms" **218** (e.g., a set of terms which, when occurring in a single communication are indicative of a particular scam story), and/or analysis based on "indicating terms" **220** (e.g., individual terms whose presence in a communication are highly indicative of scams). In some embodiments, platform **200** includes a honeypot engine **222** configured to collect fraudulent messages, along with their sender information, by generating traffic from one or more honeypot accounts; and collecting and evaluating traffic to such accounts. The indications of fraudulent activity derived from these actions can be used to help improve the filtering of messages received by real users.

In some embodiments, a tertiary classification (e.g., "bad," "good," and "undetermined") is made based on the performed evaluations (and, if applicable, taking into account any user-specific information). Where the result of the evaluation is "undetermined," the communication is optionally stored in repository **210** and provided to detailed classifier **212** for enhanced processing. In some embodiments, quick classifier **204** is provided by one entity (e.g., is located on device **102** or provided by an entity such as corporate mail service **114**), and detailed classifier **212** is provided by another entity (e.g., is provided by a third party operator of platform **200**, proxy **116**, an ISP, or other applicable entity).

In various embodiments, quick classifier **204** and detailed classifier **212** employ different classification techniques. For example, quick classifier **204** may rely solely on white/black-lists (e.g., requiring less than one second to process a message), while detailed classifier **212** may employ machine learning or other more sophisticated/resource-intensive automated review techniques (e.g., requiring two minutes of processing per message). As will be described in more detail below, in some embodiments, detailed classifier **212** makes use of one or more human reviewers (**120-124**) instead of or in addition to performing automated analysis. For example, review coordination engine **214** can make available a copy of the communication to one or more human reviewers, who determine whether the communication should be classified as "bad" or "good." The reviewer feedback is provided back to detailed classifier **212**, which uses the information to determine a final disposition/classification of the message. In some embodiments, when a message is sent out for human review, the conclusion of the human review decides the disposition of the message. In other embodiments, the human classification is treated as one factor of a score (e.g., worth 50 points), discussed in more detail below.

US 9,245,115 B1

7

In some embodiments, the reviewers are assigned reputation information (e.g., by coordinator **214**), which is stored in reputation database **216**. The reviewers can also be compensated for their reviewing efforts, with associated book-keeping being performed by coordinator **214** or another appropriate module. As will be described in more detail below, the reviewers may comprise a variety of individuals, including paid employs of the operator of platform **200**, other users of platform **200** (e.g., who perform reviews in exchange for a discount/rebate on services), a family member (e.g. Eva on behalf of Alice), and/or members of a third party outsourcing platform, such as Amazon Mechanical Turk. In some cases, such as where the human analysis is performed by a trusted entity within an organization (e.g., a member of the IT department reviewing an email sent to Bob at his work address), the full text of the message may be provided to the reviewer. In other embodiments, the message is partially redacted prior to being provided to a reviewer, also as described in more detail below.

—Tertiary Classification of Communications—

FIG. **3** depicts an example of a set of score thresholds used in an embodiment of a tertiary communication classification system. In some embodiments the set of thresholds is used for all users of a classification system (e.g., where corporate mail service **114** uses the same settings for all users). In other embodiments, the set of thresholds is adjustable on a per-user or per-user-group basis, either at the request of the user(s) or based on factors such as an assessment of the vulnerability of the user/user-group to various communication-based scams/threats.

In the example shown, a communication that receives a score (e.g., from quick classifier **204**) of less than 60 is determined to be “good.” A communication that receives a score of greater than 80 is determined to be “bad.” A communication that receives a score between those values is determined to be “undetermined” and flagged for further analysis (e.g., by detailed classifier **212**). In various embodiments, the thresholds are set such that there are no false positives: all emails for which there is a risk for false positives (i.e., a “bad” email being classified as a “good” email) are instead classified as “undetermined” and subjected to additional processing. The determination of how to set the thresholds is a risk assessment wherein the risks of false positives are weighted against the risk of false negatives.

Communications that are determined to be “good” (also referred to herein as “green”) are delivered to their intended recipient. For example, an email intended for Alice that receives a score of “10” is ultimately delivered to Alice’s inbox on web mail service **112**. The display of good messages may also be modified, e.g., so that “good” messages are colored green or include green elements when rendered.

Communications that are determined to be “bad” (also referred to herein as “red”) are not delivered, in some embodiments. One reason to not deliver the message at all, if determined to be bad, is that an unsophisticated user may unwittingly believe the message has been misclassified and fall for the scam represented by the message. Other handling of “bad” communications can also be configured. For example, “bad” messages can be delivered to a special folder, or are marked as being highly suspicious (e.g., colored bright red when displayed). In some embodiments, “bad” messages are delivered to a separate account associated with the user. As one example, a “bad” message sent by Charlie to `alice@examplewebmail.com` (Alice’s email address on service **112**) could be redirected to `alice.screened@examplewebmail.com`. Alice could authorize Eve to access the secondary account (but not her primary

8

account) to review and permanently delete any “bad” messages, and to make sure no false negatives (i.e., “good” communications erroneously classified as “bad”) occur.

As mentioned above, messages that are neither “good” nor “bad” are flagged as “undetermined” (also referred to as “yellow”) and subjected to additional processing prior to final disposition. As needed, the communication is held (e.g., in repository **210**) until a determination is made as to whether the communication is “good” or “bad.” Additional detail regarding the additional processing of “undetermined” messages is provided below.

FIG. **4** illustrates an embodiment of a process for classifying communications.

In some embodiments process **400** is performed by platform **200**. The process begins at **402** when an electronic communication is received. As one example, a communication is received at **402** when web mail service **112** (which includes at least some components of platform **200**) receives a message from Charlie addressed to Alice. As another example, where at least some of the functionality performed by platform **200** is incorporated into a mail client installed on Bob’s laptop **106**, the mail client could receive a communication at **402** when Bob’s mail client contacts corporate mail service **114** to retrieve new mail (e.g., via POP).

At **404**, the communication is classified according to a tertiary classification scheme. As explained above, in some embodiments, the communication might be definitively classified as “good” or “bad” by a quick classifier **204**. If the communication is instead determined to be “undetermined” by the quick classifier, it is provided to detailed classifier **212** for heightened review. In some embodiments, the quick classifier and the detailed classifier are collocated on a single platform (e.g., as depicted in FIG. **2**). In other embodiments, the classifiers reside on separate devices and/or may be configured by or under the control of distinct entities. As one example, a quick classifier could be included in a mail client resident on phone **104**. As the phone has limited computing and other resources, and messages received on the phone could be handled by a remote detailed classifier (e.g., provided by proxy **116**). Further, in some embodiments multiple detailed classifiers are employed, and/or multiple rounds of enhanced scrutiny are applied to messages that are not clearly “good” or “bad.” As one example, where detailed classifier cannot definitively determine whether a message is “good” or “bad,” the message can be provided to one or more amateur human reviewers (e.g., members of the public who have agreed to help review messages). If the amateur reviewers are similarly unable to determine/agree that a given message is “good” or “bad,” the message can be provided to professional reviewers (e.g., employees of the owner of platform **200** or paid contractors). Protected users/subscribers can also potentially act as reviewers (whether for themselves or others). For example, if a sufficient number of protected users report a received message as “spam,” the message would be considered “bad,” and/or would cause the message to be provided to be given to amateur or professional reviewers to classify.

The classification performed at **404** can be based on a numerical score (e.g., using numerical thresholds such as are depicted in FIG. **3**). The classification performed at **404** can also be based on a set of rules. Examples of both approaches to classification are given below, and will draw from the following list of example considerations:

1. “collection of terms”—The co-occurrence of certain terms from separate domains in a message can be indicative of a fraudulent message that corresponds to a particular scam scenario (described in more detail below). As one example, a message that contains the term “MCTN” (a term specific to

US 9,245,115 B1

9

Western Union) and also the term “Google Wallet” is indicative of fraud. Scammers frequently offer to send fake Google Wallet payments and request money back using Western Union. The two terms are extraordinarily unlikely to co-occur in a legitimate email discussion. However, a term such as “Google Wallet,” by itself, could be prevalent in legitimate emails; a blanket blacklisting of the term is likely to result in far too many false positives (flagging legitimate messages as scam messages) to be tolerated by users being protected. The presence of a collection of terms in a message almost certainly indicates the message is fraudulent. Another example collection of terms is: “Internet Lottery,” “your email has won,” “congratulations,” and “million dollars.” The last term, “million dollars” is also considered present in a message if any so-called equivalent terms are present; such terms may consist of a list “millions dollars”, “million pounds”, and “several millions.”

2. “indicating terms”—Terms that are statistically common in scam communications and uncommon in legitimate communications. “Internet” and “your name” are not indicating terms, as they are very prevalent in legitimate communications. “Abacha,” however, is virtually absent from legitimate communications but prevalent in scam communications. Additional examples of “indicating terms” include “modalities,” “no risk,” “(\$*,000,000)” where * denotes an arbitrary value. The absence of any indicating terms in a message almost certainly indicates that the message is benign.

3. “friendly email”—A user-specific whitelist of email addresses. In some embodiments, the user or an administrator provides the whitelist. In other embodiments, the whitelist is constructed based on an observation (e.g., by platform 200) of communications. As one example, once a protected user has exchanged more than a threshold number of emails with someone at a given email address, over a threshold period of time (e.g., ten emails over the course of two weeks), then the address could be designated as a friendly email.

4. “friendly location”—A user-specific geolocation of friendly emails (or other identifier of a communication’s origin). In some embodiments, the geolocations associated with email addresses that are considered to be friendly (e.g., per above) are designated as “friendly” after threshold requirements are met. As one example, if Alice has a friend in Paris, France, with whom she routinely communicates, that friend would eventually be added to the “friendly email list.” The geolocation of “Paris, France” could then be added as a friendly location (e.g., after Alice has exchanged more than twenty messages with her Parisian friend, over a period of at least one month). In some embodiments, a freshness test is employed, so that a location only remains friendly so long as the threshold amount of communication continues to be exchanged (or some other “maintenance” level of communication). An example of a way the likely approximate location of the sender can be determined is by review of the routing path, which is always available to a receiving ISP, and commonly available to the final recipient.

The geolocations designated as “friendly” can take a variety of forms, ranging from specific IP blocks/subnets (e.g., indicative of the particular French ISP used by Alice’s friend), to political boundaries such as neighborhoods/ZIP codes/cities/counties/states/countries, or arbitrary designations such as “within a 20 mile radius” of a given location. The classification can also be performed by determining if the IP is not in a given range, e.g., “any not originating in Denmark” is not friendly.

Alice might visit her friend (and make new friends) or otherwise legitimately begin communicating with others in

10

Paris. Such communications are much less likely to be fraudulent/take advantage of Alice than communications originating from a geolocation with which she’s not previously had contact (e.g., Spain or Ghana). A variety of rules can be used to govern whether/when a friendly location is added. For example, if Alice’s Parisian friend visits China for a week, and sends a handful of emails to Alice, platform 200 will not automatically add the country of China to Alice’s friendly location list. One way to ensure China is not added is to require a threshold number of additional “friendly email” addresses in a given region before adding a region, and/or connecting the number of friendly emails to the footprint of the associated geolocation (e.g., requiring Alice to receive mail on a regular basis from three people in a given state or country before adding the entire state/country).

In some embodiments, friendly languages are determined, and may, for example, correspond to the languages spoken in the friendly locations. Thus, if a language used in the message is either on a particular list of languages, or is not on a list of on a list of particular languages, then this can be used as a factor to cause the email to be identified as good, suspect, or bad. Rules can also be combined, describing scenarios such as: “All email from outside Scandinavia is considered suspect, except email from Britain if the language is Swedish and there are no indicating terms in the email.” A variety of approaches can be used to assess messages according to such rules. For example, services such as Google translate can be used; comparisons of message contents to word lists of different languages can be performed; and/or the encoding of the message and the associated language character table can be identified. Most computers use UTF (8,16) to display contents. In the case of HTML pages, the page typically has a metatag that indicates the encoding of the page, and if the characters are from a certain part of the coding table that shows the language of the page.

While it may seem unnecessarily restrictive to block traffic from entire subnets, countries or even continents, or likewise to block traffic in certain languages, there are many users to whom the Internet poses greater threats than benefits unless drastic limitations of freedom are made. Thus, to such users, or to concerned family members, it may be desirable/reasonable to block traffic from all countries where the protected user does not have any friends, family or business partners, or conversely, to only admit traffic from whitelisted locations. The importance of a given factor, including friendly location, can be determined by weights associated with the rules; also, for some users, only a subset of the rules need be active or configured.

5. “suspect location”—A listing of VPN proxy addresses, Tor exit nodes, zombie/bot nodes, and other known-bad sending locations that is not user-specific. As one example, if a particular ISP in Ghana is known for originating a great deal of scam messages, that ISP could be designated as a “suspect location.” And, paths can also be used, e.g., any web email originating in Ghana, or originating in Spain.

6. “undisclosed location”—A listing of webmail providers that is not user-specific, and a location that resolves to a VPN, known bot node, or similar problematic location.

7. “global friendly senders”—A listing of well-known, benign electronic commerce and other service providers that is not user-specific.

8. “spoof”—Messages sent to “undisclosed recipients” and/or other indicators of sender spoofing. Additional examples include: (1) comparing the originating IP address with the domain; (2) identifying suspect IP addresses on the path of the message; (3) identifying an unusual number of hops; (4) identifying previously identified bad IP addresses in

US 9,245,115 B1

11

long headers; (5) email contents being misaligned with the domain appearing to have originated the message; (6) email contents being misaligned with the IP addresses on the path of the message; and/or (7) the email has a sufficiently different reply-to address from the apparent sender address, or (8) the email has a sufficiently different reply-address from both the apparent sender address and the originating domain.

Example

Numerical Scoring

Each of the above eight example considerations is associated with a number of points. The following is one example of how points could be assigned:

(collection of terms fires): 65 points

(indicating terms fires): 10 points

not (friendly email): 25 points

not (friendly location): 25 points

(suspect location): 30 points

(undisclosed location): 10 points

(global friendly sender): -20 points (this value is negative, indicating that the presence of the condition being true is indicative of a "good" message)

(spoof): 50 points

For a given communication, the points are summed, and compared to thresholds. Below are two examples of thresholds, one set for Alice (e.g., by Eve) and one set for Bob (e.g., by his employer):

Alice:

<15 points—green

>50 points—red

otherwise yellow

Bob:

<25 points—green

>60 points—red

otherwise yellow

Example

Rule-Based

Various rules involving the eight example considerations can be defined and associated with tertiary classifications. The following are two examples of how rules can be specified—again, with Alice's rules being more strict against potential scam, and Bob's being more lax:

Alice:

RED if:

(Collection of terms fires), or

Spoof, or

no color determined and (not friendly location) and (not friendly email) and (not global friendly senders) and (indicating terms), or

no color determined and ((undisclosed location) or (suspect location)) and (indicating terms fires)

GREEN if:

no color determined and (friendly email) or (global friendly senders), or

no color determined and (friendly location) and not (indicating terms fires)

YELLOW otherwise.

Bob:

RED if:

(Collection of terms fires), or

Spoof

GREEN if:

12

no color determined and (friendly email) or (global friendly senders), or

no color determined and (friendly location) and not (indicating terms fires), or

no color determined and (friendly location), or

no color determined and (friendly location) and not (indicating terms fires)

YELLOW if:

no color determined and (not friendly location) and (not friendly email) and (not global friendly senders) and (indicating terms), or

no color determined and ((undisclosed location) or (suspect location)) and (indicating terms fires), or

[otherwise doesn't match any rules].

The rules for what is red, green, and yellow are configurable on a per-case basis and an administrator can select and configure these rules. Conflict between rules can be managed by a pessimistic approach (if any indicator says it is red, then it is red); using a threshold (if at least X indicators say it is red, then it is red); or with exceptions (it is not red if it is on the whitelist, otherwise if any indicator says it is bad then it is red.) These different approaches carry different risks of false positives, where the risk for false positives would be higher for the pessimistic approach than it would be for the other described approaches. Yet other more flexible policies for determining how to manage conflicts can also be used; such as by having each rule associate a score to each message and each rule being associated with a weight, allowing an aggregated weight to be computed and compared to a threshold value, which can be a scalar or a vector element.

FIG. 5A illustrates an example of an electronic communication. In particular, message 500 is an example of a fraudulent email message that Charlie has attempted to send to Alice. In this example, Charlie has previously contacted Alice (e.g., letting her know she has won a prize and asking for her address to determine which prize she has won), and Alice has responded (e.g., with her address). In practice, Charlie's original message, and/or the concatenation of the exchanges between Charlie and Alice would have been classified by quick classifier 204 as "bad" (e.g., based on the "collection of terms" technique described below). However, for purposes of this example, suppose that quick classifier 204 has classified message 500 as "undetermined." The message is provided to detailed classifier 212 for further analysis. As mentioned above, in some embodiments, detailed classifier 212 is configured to leverage the assistance of human reviewers in determining whether an undetermined message should be definitively classified as good or bad. Detailed classifier 212 provides the message to review coordinator 214 which redacts the message as applicable, and provides the redacted message to one or more human reviewers. In particular, personally identifiable information, such as contact information associated with the protected user (i.e., potential victim recipient) is redacted. In some embodiments, other processing is also performed prior to providing the message to a reviewer for review. For example, terms commonly used by scammers can be highlighted, and path information made easier to understand, as described in more detail below. Such processing can be performed based on parsing of the text (or optical character recognition (OCR) of images, followed by parsing of the resulting text); comparisons to known personally identifiable information (PII) terms/formats (e.g., common first names; recognition of phone numbers; recognition of addresses); and comparison to known terms commonly used by fraudsters, but not as commonly used in non-fraudulent messages (described in more detail below). In some scams, messages are included as attachments, and/or links

US 9,245,115 B1

13

included in the email (or, e.g., an SMS) direct victims to a website that includes an image of text. Processing can further include examining attachments, and detecting and following such links, and OCR'ing/parsing the obtained content as applicable.

FIG. 5B illustrates an example of an interface for classifying an electronic communication. The example shown is an embodiment of an interface shown to a reviewer, such as reviewer 120. The interface can be provided in a variety of ways. As one example, platform 200 may provide a web interface/portal which reviewers can access, log into, and then select an option to "start reviewing." As another example, e.g., where Eve is reviewing Alice's messages, Eve may receive an email or other alert, letting her know that new mail which requires review has arrived for Alice, and asking Eve to access an interface provided by platform 200. In yet other embodiments, the interface is provided as a tablet or other mobile device app, allowing reviewers to review messages in a custom interface.

In interface 550, Alice's email address 502 has been redacted (552). The sender's email address 504 is underlined (554) to indicate that the sender is involved in other messages, which the reviewer can access by clicking on region 554. Subject line 506 is shown in the redacted version of the message (556). The path of the message 508 is textually and/or visually clarified/simplified. For example, the IP address information is replaced with a geographic location and an indication that it may be spoofed (558). Other techniques for simplifying location information can also be provided, such as by showing country information on a map when the reviewer hovers a mouse pointer over region 558.

Alice's name 510 is identified as a personal name and removed (560), as is component 512 (562), which is identified as an address. In some embodiments, instead of blacking out the personal information, the information is placed with an indicator of what has been removed, e.g. "NAME" in region 560 and "ADDRESS" in region 562. Where reviewer coordinator 214 is unable to definitively determine whether a portion of the message should be redacted, the label over the redacted portion can indicate as such, e.g., "PROBABLY AN ADDRESS." Elements of the remaining text commonly associated with scams are highlighted for the reviewer's convenience (564-568).

In region 570, the reviewer is asked to make a determination of how the message should be classified, by clicking on one of buttons 572, 574, or 576. The result of a button press is received by review coordinator 214, which collects the feedback from any additional reviewers who have been asked to review the message. In some embodiments, a feedback field or other mechanism is included in the interface so that the reviewer can provide feedback on why the message was classified. As one example, the reviewer could explain what reasons led the reviewer to be "not sure" about the message, to help a subsequent reviewer come to a definitive decision.

In some embodiments, the results of other filter rules are indicated in the image shown to the reviewer, to assist the reviewer in assessing the message. For example, where the sender location is friendly, that information could be provided at the top of the interface, or inline (e.g., next to region 554). As another example, information such as "UK proxy" or "unknown ISP" can be included in region 558.

In some embodiments, a single definitive (i.e., "good" or "bad" designation, but not "not sure" designation) classification by a reviewer is sufficient to classify the message. This may be the case, for example, where a relative of the protected user is performing the review (i.e., Eve reviewing on behalf of Alice), or a designated employee is performing the review

14

(i.e., a member of Bob's company's IT department reviewing on behalf of Bob). In this scenario, button 574 may be omitted, if applicable. The number of reviewers assigned to review a message can also be based on a service level associated with the protection. For example, platform 200 may offer free protection services (where only one human reviewer will resolve undetermined messages) and also offer premium services (where multiple reviewers will vote, and/or where the experience of the reviewers varies based on subscription level).

Where multiple reviewers are asked to review a message (e.g., in parallel, as a group), if a sufficient number of reviewers indicate that a given message is fraudulent by clicking on the "bad" button 576, then the message is classified as "bad." If a sufficiently large number of reviewers select the "good" button 572, then message is considered good. If a sufficient number reviewers select option "not sure" 574, in some embodiments, the classification task is elevated to more experienced reviewers, e.g., as assessed by the number of cases they have judged, their reputation, and/or the duration that they have provided feedback. A variety of approaches can be used to determine whether the "sufficient" number is reached. As one example, a rule can be specified that the selected group of reviewers must agree unanimously. As another example, a rule can be specified that a threshold percentage of the group must agree. As yet another example, the "votes" of the reviewers can be weighted by their respective reputation scores (e.g., stored in database 216). Additional information regarding reviewer reputations is provided below.

Reputations can be assigned to reviewers in a variety of ways. As one example, reviewers can be asked to review training materials and then, upon completion, evaluate a set of sample messages. A reviewer's reviewing accuracy with respect to the sample set can be used to assign a reputation to the reviewer. In some embodiments, the reputation is binary (e.g., "trusted" or "not trusted") or tertiary (e.g., "good reviewer," "average reviewer," "novice reviewer"). The reputation can also be a score (e.g., 16/20) or percentile (e.g., 75th percentile). In some embodiments, novice reviewers are assigned a starting reputation of 10/100. As they classify messages, their score is increased or decreased based on whether other, more experienced reviewers agree with their assessment. If a novice reviewer's score reaches a low threshold (e.g., 5/100), the reviewer is warned to do a better job, and/or prevented from reviewing any more messages, due to poor performance (e.g., once the reviewer's score dips below 3/100). Where compensation is awarded (whether in the form of monetary payment, or additional reputation/other points), the compensation for a "correct" vote may be evenly distributed across all reviewers voting correctly, and may also be distributed according to a formula, e.g., that takes into account the reviewer's respective reputation scores, history, etc.

FIG. 6 depicts an example of a review performed by multiple reviewers. In the example shown, reviewer A (602) has a reputation of 46. Reviewer B (604) has a reputation of 72. Reviewer C (606) has a reputation of 51. Reviewers A and B have both flagged the message being reviewed as "unsafe" (i.e., bad). Reviewer C has flagged the message a "safe" (i.e., good). One way to determine an assessment of the message is to sum each of the votes. In the example shown in FIG. 6, such a score could be computed as 46+72-51 (total 67). Another way to determine the assessment is to assign one vote to each reviewer, and then weight the respective votes based on reputation buckets. As one example, a reputation at or above 65 could be accorded a full vote, a reputation between 50 and 65 could be accorded 0.75 votes, and a reputation 50 or below

US 9,245,115 B1

15

could be accorded 0 votes. Votes determined to be correct are rewarded with heightened reputations, and votes determined to be incorrect are penalized with lowered reputations. In some embodiments, a function is used to determine how much reward or penalty is to be applied. As one example, reputations may be recomputed daily or weekly, and the total number of messages reviewed by a reviewer taken into account considered when adjusting the reputation. As one example, a reviewer who reviews at least 20 messages in a week, and gets no more than ten percent wrong, could be assigned two points of reputation for that week. A reviewer who gets more than ten percent wrong (irrespective of total number reviewed) could be penalized by four points.

In various embodiments, reviewers are made aware of the fact that there are at least some additional reviewers reviewing the message. Information about the other reviewers, such as their number and respective reputations, can selectively be displayed or not displayed (e.g., based on configuration settings made by an administrator). In some embodiments, reviewers are unaware (or are otherwise not informed) about whether or not multiple reviewers are examining the message they are examining

—Collection of Terms—

Overview

As mentioned above, in various embodiments, platform 200 evaluates communications based on the presence of a “collection of terms.” As explained above, the presence of a collection of terms in a message almost certainly indicates the message is fraudulent, and the absence of any indicating terms in a message almost certainly indicates that the message is benign. Accordingly, in some embodiments, evaluation of a communication is performed using a collection of terms, and a binary classification (i.e., scam or not scam), score, or other non-tertiary classification scheme is used. In other embodiments, the presence of a collection of terms in a communication (and/or absence of indicating terms) is one factor in tertiary classification.

In some embodiments, quick classifier 204 and detailed classifier 212 cooperate to perform tertiary classification of messages. In other embodiments, at least some portion of the time, the classifiers use a “collection of terms” to classify communications (e.g., as fraudulent or not) and do not perform tertiary classification. In various embodiments, a single classifier is included in platform 200 (e.g., quick classifier 204 is omitted, and detailed classifier performs any functionality otherwise provided by quick classifier 204).

FIG. 7A depicts an example of a collection of terms. In particular, FIG. 7A provides an outline of a particular form of scam that is perpetrated by 419 scammers—trying to convince the victim that he or she is entitled to a large sum of money as an inheritance, and that the money will be provided as soon as the victim pays a small fee. Although the main points of the message will be common across all such scam messages conforming to the “story” the scammer is trying to trick the victim into believing, the actual wording of the scam message may vary from message to message (e.g., to thwart detection, because the message has a particular author with a distinctive writing style, or because the message was written in a first language and translated to a second). Further, subtle variations may occur due to writing problems such as misspellings.

Each row in the collection of terms depicted in FIG. 7A corresponds to one aspect of the inheritance scam story. Where multiple terms appear on a given row, the terms are collectively referred to as an equivalence class—terms that fulfill the same purpose if used in the story. For example, the particular scam represented by FIG. 7A typically begins with

16

an introduction of either “My name is” (702) or “I am” (704). The scam will next invoke a long-lost relative (or their representative). Equivalence class terms for this aspect of the story are shown in region 706. Next, the scam will describe the large amount of money (in one of three formats shown in region 708) that can be collected by the victim in one of three formats. The scam then indicates that all that is required for the victim to receive the money (e.g., “transfer” 710) is for the victim to provide banking details (see region 712 for terms). The victim is encouraged to provide the banking details right away (see region 714 for terms), e.g., to minimize the likelihood the victim will tell a friend or relative about the email and be discouraged from providing payment information.

FIG. 7B shows an example of a fraudulent message that would be detected based on analysis by content evaluation engine 208 of the collection of terms depicted in FIG. 7A. The terms in message 750 that are present in the collection of terms of FIG. 7A are underlined. In some embodiments, which term in an equivalence class is used in a message (e.g., “My name is” vs. “I am”) is not taken into account when evaluating the message. In other embodiments, different terms receive different scores. As one example, “huge sum” might be scored higher (i.e., indicating the message is more likely to be fraudulent) than “;000.”

FIG. 8 illustrates an example of a process for classifying communications. In some embodiments, process 800 is performed by platform 200. The process begins at 802 when an electronic communication is received. As one example, a communication is received at 802 when web mail service 112 (which includes at least some components of platform 200) receives a message from Charlie addressed to Alice. As another example, where at least some of the functionality performed by platform 200 is incorporated into a mail client installed on Bob’s laptop 106, the mail client could receive a communication at 802 when Bob’s mail client contacts corporate mail service 114 to retrieve new mail (e.g., via IMAP).

At 804, the communication is classified using a collection of terms. As explained above, in some embodiments the communication might be definitively classified as “good” or “bad” based on the analysis of the message against the set of collections of terms 218. In other embodiments, the collections of terms analysis is one consideration among multiple considerations (e.g., the additional example considerations listed above). In various embodiments, the distance between at least some terms appearing in the message is taken into account when determining whether the message should be marked as fraudulent based on the presence in the message of a collection of terms. As one example, while presence of the terms, “Nigeria” and “senator” in the same short message may typically indicate that the message is fraudulent, the message is likely not fraudulent where the terms are separated by 5,000 characters.

The classification performed at 804 can be performed using a variety of techniques. For example, a collection of terms can be evaluated using a rule-based approach (e.g., testing for the presence of words, and/or applying a threshold number of words whose presence are needed for a match to be found); using a support vector machine, where the elements of the support vector corresponds to terms or words; and/or using general artificial intelligence methods, such as neural networks, wherein nodes correspond to terms or words, and wherein the values associated with connectors cause an output corresponding essentially to a rule-based method. In each of the aforementioned embodiments, a value associated with the severity of the collection of terms being identified can be generated and output, where multiple values are generated if multiple collections of terms have been identified.

US 9,245,115 B1

17

Additional Information Regarding Collections of Terms

In some embodiments, each term (or its equivalent) must appear in the message in the order it appears in the collection. Thus, using the example of FIG. 7A, in some embodiments, if “transfer” appears before “huge sum” in a message being analyzed, the message will not be flagged as a scam, because the ordering in the collection of terms is reversed. In other embodiments, order of terms does not matter, e.g., so long as at least one term from each line of the collection shown in FIG. 7A is present in the message, the message will be classified as an inheritance scam.

In some embodiments, platform 200 maintains scores associated with each collection of terms. One such value indicates, for each type of scam, how successful the associated term collection is at matching fraudulent emails making use of that scam. Based on factors such as the concern for various types of scams, and based on computational limitations, a selection of which term collections are to be used can be made, e.g., where processing is performed on a device with limited resources, such as phone 104.

A second value associated with each collection of terms indicates the risk for false positives associated with the term collection, in the context of a given user. Example ways to determine the value is by scanning the user’s inbox; by letting the user identify his or her normal activities; and/or by running the system for some amount of time; and determining the value based on classification of uncertain cases by human reviewers who review messages and classify them. This second value can also be used to select collections of terms, e.g., to avoid term collections that lead to higher false positive rates than a particular user find acceptable.

Both values can be configured based on the preferences of the protected user, and on the service level of the user (e.g., where users with higher service levels are given higher computational effort). In some embodiments, a collection of terms is matched to a portion of an email address, and a determination is made as to whether the email is from a domain associated with the terms; if it is not, then the email is flagged. As one example, an email with terms suggesting that the email is the confirmation of a financial institution payment but which is not sent from the financial institution domain is flagged as scam. In another example, a determination is made as to whether the message is from a particular sender, and if it is not, then the message is flagged as scam. In yet another example, all words are normalized before the comparison is made. This includes performing a consistent capitalization, correcting likely spelling mistakes by replacing words with the most likely candidates from a list of related words, where this list is created to emphasize words commonly used by scammers.

The following is another example of detecting a fraudulent message using a collection of terms. Suppose there are a total of two terms included in the collection (corresponding to a fraud in which victims are asked to send money by Western Union in exchange for a bogus Amazon.com payment). In this example, no equivalence terms are included—just a total of two distinct terms—(“Western Union”, “Amazon payment”). If a document contains both of these terms, whether separated by other words or not, then the document is considered to match. Suppose the message is, “Here is an Amazon payment for \$100. Please send me \$50 with Western Union.” Such a message would match the collection of terms, as would “Please send your Western Union payment after you receive the Amazon payment.” However, a message of, “Here is an Amazon payment for the Western Digital hard drive I want to

18

purchase. Please send it to my home in Union, N.J.,” would not match since “Western” and “Union” are separated. A message of, “Here is an AMAZON payment for \$100, please send the money with western union” would match, where normalization is applied to remove capitalization. In an embodiment where spelling errors are corrected/normalized, “Here is an Amaz0n payment. Please send money using western union,” would match the collection of terms, since “Amaz0n” once corrected would become “Amazon,” and “unjon” would be corrected to “union” before the verification is made.

In some embodiments, a global list of equivalent terms is maintained (e.g., usable across multiple collections of terms), such as “USD,” “us\$,” and “euro.” While a Euro is not the same as a USD, the usage of either concept by a scammer is functionally the same. In some embodiments, as a message is evaluated (e.g., by content evaluation engine 208), it is first normalized by capitalization and spelling normalization, then the system replaces any terms found in the document matching a term in the list of equivalent terms with a representative term, such as the first term in the equivalence list. After that, the document is verified to determine if it matches any of the rules, such as the (“Amazon”, “Western Union”) rule. In some embodiments, any images included in/attached to/linked to in the message, are interpreted using OCR techniques, and any associated texts combined with ASCII text material before the verification is made.

In some embodiments, each of the non-equivalent terms in a collection of terms (e.g., “long lost” and “huge sum”) are associated with one or more pointers, and ordered alphabetically. The number of pointers associated with each term is the same as the number of rules for which that term is used. Each rule is represented as a vector of Boolean values, where the vector has the same length as the associated rule contains words. All the binary values are set to false before a message is parsed. The message is parsed by reviewing word by word, starting with the first word. If the word being reviewed does not fully or partially match any of the alphabetically ordered terms, then the next word is reviewed instead. If a word matches a term fully, then all Boolean values that are pointed to by the pointers associated with the term that the word matches are set to true. If one or more words matches a term partially by being the first words in the term, then the next word of the message is being added to the comparison and it is determined whether the previously partially matching words now partially or fully match any of the terms that was previously partially matched. If a full match is achieved, then the Boolean values associated with the pointers of this term are set to true. If a partial match is achieved, then the next word is added, and the process repeated. If a sequence of words being matched first partially matches and then does not match, then the system again will consider one word, starting with the second word of the previous sequence. After the entire document has been parsed in this manner, the system determines whether any of the vectors of Boolean values is all true, and if this is so, then the algorithm outputs that there is a match; otherwise it outputs that there is no match. A match means that the message is dangerous. This comparison can also be made each time a Boolean value is set to true by determining if the vector in which this Boolean value is an element is all true, and if so, output “match” and conclude the processing of the message. In a variant implementation, the system determines how many of the vectors are set to all-true; and outputs a counter corresponding to this number. Alternatively, each vector is associated with a weight, and the system determines the sum of all the weights for which the associated vectors are all-true. The message is then identified as having

US 9,245,115 B1

19

dangerous content, and the sum determines the extent of the danger. In one embodiment, the Boolean vectors are not set to all-false between the scan of two related messages that are part of a thread and sent to the same person. This provides detection capabilities in situations where information is dispersed over multiple related messages, which causes the thread of messages to be considered dangerous.

—Obtaining Collections of Terms—

Collections of terms **218**, an example of which is depicted in FIG. 7A, can be included in platform **200** in a variety of ways. As one example, a human administrator (or contractor linguist, or other appropriate entity) can manually create a given collection (and optionally assign it a title, as applicable, such as “inheritance scam”), which can be stored for use by platform **200**. As another example, messages that are flagged (e.g., by human reviewers) as being fraudulent, but are not otherwise flagged by platform **200** can be examined—either automatically, or in cooperation with humans, such as an administrator or reviewers, and collections of terms formulated to identify such fraudulent messages in the future.

FIG. 9 illustrates an example of an interface configured to receive feedback usable to create collections of terms. In the example shown, an administrator is reviewing feedback provided by three reviewers about why a particular message is believed to be fraudulent. Specifically, while interacting with an interface such as a modified version of interface **5B**, reviewers were asked to indicate which terms they believed were most important in reaching their determination of bad, by highlighting the terms prior to clicking “bad” button **576**.

The terms selected by each of the three reviewers are indicated to the administrator as three types of boxes—solid boxes indicate a selection by a first reviewer; dashed boxes indicate a selection by a second reviewer; and dotted boxes indicate a selection by a third reviewer. In the example shown in FIG. 9, the administrator is not authorized to see the full message, so certain terms (e.g., term **902**) are redacted, even for the administrator. All three reviewers agree that term **904** is probative of why the message is fraudulent. Other terms have votes from only two (e.g., **906**) or just one (e.g., **908**) of the reviewers. In various embodiments, the administrator can review the selections made by the reviewers, and act, e.g., as a fourth reviewer, to pick which terms should be included in a collection of terms usable to detect the scam represented by the message. The administrator can also set thresholds (e.g., minimum of two votes needed, reviewer reputation score needed, etc.) for automatically selecting terms, and then retain the ability to approve or veto the automatic inclusion of the collection of terms in collection **218**. In some embodiments, the flagging of terms in the message is presented to users as a CAPTCHA.

In some embodiments, automated techniques are used to generate collections of terms (and/or indicating terms). For example, suppose the classification of a given message is “bad.” Platform **200** can be configured to identify terms that distinguish it from messages of the good message set, using the TF-IDF (term frequency inverse document frequency) principle. A limited number of such terms are selected, where the number is either a system parameter or a function of the TF-IDF value, and where the terms are selected in order of decreasing TF-IDF values; while selecting at least a threshold number of word terms; at least a threshold number of bigrams; and at least a threshold number of trigrams. These selected terms are stored, and referred to as temporary terms. Platform **200** then computes a modified TF-IDF value for the normalized message and messages of the good message set, using constellations of the temporary terms, where a constellation is an unordered list of elements selected from the tem-

20

porary terms, for different such selections. This identifies collections of elements from the set of temporary terms that are particularly rare in good messages. A threshold number of the resulting terms are kept, selected in order of decreasing modified TF-IDF value. The threshold is either a parameter number or a function of the modified TF-IDF number. The result are rules that identifies the input message as bad, and the inverse of the modified TF-IDF number is an estimate of the false positive rate for classification of messages using the associated rule. These rules are then ordered in terms of decreasing values of a counter measuring how many messages in the collection of known bad messages that each such rule matches. These counters are estimates of how general the associated rule is. One or more rules are selected from the rules, where the selection criteria are low false positive rates and large degree of generality. An example selection picks the rule that maximizes a measure equaling the generality measure divided by the false positive rate, i.e., the associated counter times the associated modified TF-IDF value. The selected rules are added to the database of rules. This approach is used to compute new rules to identify bad messages. In one version of the algorithm, the entire set of known good messages is used in place of the at least one message that is part of the input.

As another example, collections of terms can be generated using artificial intelligence techniques configured to identify common words in scam messages, but which are not as common in desirable messages; identify collections of such words that are frequent in scam messages but which are highly infrequent in desirable messages; and identify collections of such terms that are common in scam messages but which are essentially absent in desirable messages.

—Temporal Considerations—

The disclosed techniques can take into consideration temporal relationships between messages when making an assessment. For example, platform **200** can be configured to scan sequences of messages forming a conversation. It may be that one of the messages in the sequence does not have sufficient evidence of being abusive, whereas a sequence of such messages collectively provides sufficient evidence to be filtered out as being bad. This will cause any future emails of the same type or in the same sequence to also be considered bad.

FIG. 10 illustrates an example of such a sequence of messages. In the first message (**1002**), a user called “Grandma” receives a seemingly benign email from someone claiming to be a long lost friend. It does not mention lotteries. Grandma responds (**1004**) that she cannot remember her friend, then gets a second email (**1006**) saying that they were in the same elementary school, and now her friend sells lottery tickets and has five grandchildren. Grandma responds (**1008**) that this sounds like a fun thing to do, and that she has ten grandchildren. Her “long lost friend” then says (**1110**) that the reason she contacted Grandma was that she saw her name as one of the lottery winners, and remembered her name from her childhood, then decided to find her to tell her about her winnings. How could she not pick up the money, it is nearly a million dollars, and all she has to do is to pay the processing fee of \$565.

Each email in the exchange, by itself, might be seen as innocuous, with the potential exception message **1110**. By the time message **1110** is received, however, most existing spam filters would have whitelisted the scammer, given the number of emails sent and received from her by Grandma without incident. In various embodiments, platform **200** examines the entire sequence of emails (or a moving window of several emails), concatenating the text together and performing

US 9,245,115 B1

21

analysis on the concatenated text. The concatenated text would readily match a “Lottery Scam” collection of words, and the messages would be classified as “bad,” accordingly.

A second example of temporal processing is as follows. Suppose a user is receiving a sequence of emails over a few weeks time, where the sequence of emails establishes an online friendship or relationship, and then asks for money for some purpose. The initial sequence of emails is purely intended to establish trust, after which the typical request for money arrives. A person who has seen such a scam perpetrated might recognize its making from the early emails. A machine learning component (e.g., of content evaluation engine 208) can identify a sequence of messages as bad when identifying the request for money, and then identify indications in the trust-establishing emails that are indicative—whether by themselves or as a subsequence—of the request to come. This way, the machine learning component will constitute an early-warning system in which indications of fraud are picked up before there are signs that by themselves correspond to an effort to extract money.

Although the foregoing embodiments have been described in some detail for purposes of clarity of understanding, the invention is not limited to the details provided. There are many alternative ways of implementing the invention. The disclosed embodiments are illustrative and not restrictive.

What is claimed is:

1. A system, comprising:
a processor configured to:
receive an electronic communication;
obtain, from a data store, a first collection of terms,
wherein:
the first collection of terms corresponds to a first attack strategy;
the first attack strategy is associated with a plurality of narrative points, wherein the plurality of narrative points comprises a plurality of equivalence classes;
each narrative point of the first attack strategy is associated with a corresponding group of terms representative of the narrative point;
the groups of terms corresponding to the narrative points collectively comprise the first collection of terms; and
a second attack strategy that is different from the first attack strategy is associated with a second collection of terms;
evaluate the received communication to determine whether the received communication is indicative of the first attack strategy, wherein the evaluating includes:
for each narrative point of the first attack strategy, determining whether the received communication includes the presence of one or more terms matching to the group of terms corresponding to the narrative point;
determining whether the received communication includes the presence of terms matching to a threshold number of narrative points; and
determining whether terms present in the communication that match to the narrative points of the first attack strategy occur in a particular order; and
classify the received communication based at least in part on the evaluation; and a memory coupled to the processor and configured to provide the processor with instructions.
2. The system of claim 1 wherein the processor is further configured to perform a language detection on the received communication.
3. The system of claim 2 wherein the language detection is performed based at least in part on a detection of the presence of a set of terms.

22

4. The system of claim 1 wherein the processor is further configured to determine that the received electronic communication is associated with at least one additional message.

5. The method of claim 4 wherein the processor is further configured to classify the received communication at least in part by examining the received communication and the additional message in the aggregate.

6. The system of claim 1 wherein classifying includes determining a distance between terms as they appear in the communication.

7. The system of claim 1 wherein the first collection of terms is determined in response to an evaluation by a plurality of workers of a fraudulent communication.

8. The system of claim 1 wherein the first collection of terms is determined based at least in part on a prior classification by a server of a message as fraudulent.

9. The system of claim 1 wherein the first collection of terms includes, for at least one term in the collection, a set of equivalent terms.

10. The system of claim 9 wherein the processor is substitute a term in the set of equivalent terms for a term appearing in the communication.

11. The system of claim 1 wherein a term includes at least one letter and at least one number.

12. The system of claim 1 wherein the processor is configured to analyze a portion of an email address associated with the communication and to select a collection of terms for use in classification based at least in part on the analyzed address.

13. The system of claim 1 wherein the communication is received from a honeypot.

14. The system of claim 1 wherein the communication is classified in accordance with a tertiary classification scheme.

15. The system of claim 1 wherein the processor is further configured to receive, from at least one human reviewer, a classification of the communication.

16. The system of claim 1 wherein the classifying is performed by using a set of rules.

17. The system of claim 1 wherein the classifying is performed at least in part by using a support vector machine.

18. The system of claim 1 wherein the classifying is performed at least in part by using an artificial intelligence computation.

19. A method of, comprising:
receiving an electronic communication;
obtaining, from a data store, a first collection of terms,
wherein:

the first collection of terms corresponds to a first attack strategy;

the first attack strategy is associated with a plurality of narrative points,

wherein the plurality of narrative points comprises a plurality of equivalence classes;

each narrative point of the first attack strategy is associated with a corresponding group of terms representative of the narrative point;

the groups of terms corresponding to the narrative points collectively comprise the first collection of terms; and

a second attack strategy that is different from the first attack strategy is associated with a second collection of terms;
evaluating, using a processor, the received communication to determine whether the received communication is indicative of the first attack strategy, wherein the evaluating includes:

for each narrative point of the first attack strategy, determining whether the received communication includes the presence of one or more terms matching to the group of terms corresponding to the narrative point;

US 9,245,115 B1

23

determining whether the received communication includes the presence of terms matching to a threshold number of narrative points; and

determining whether terms present in the communication that match to the narrative points of the first attack strategy occur in a particular order; and

classifying the received communication based at least in part on the evaluation.

20. A computer program product embodied in a non-transitory tangible computer readable storage medium and comprising computer instructions for:

receiving an electronic communication;

obtaining, from a data store, a first collection of terms, wherein:

the first collection of terms corresponds to a first attack strategy;

the first attack strategy is associated with a plurality of narrative points,

wherein the plurality of narrative points comprises a plurality of equivalence classes;

each narrative point of the first attack strategy is associated with a corresponding group of terms representative of the narrative point;

24

the groups of terms corresponding to the narrative points collectively comprise the first collection of terms; and a second attack strategy that is different from the first attack strategy is associated with a second collection of terms;

evaluating the received communication to determine whether the received communication is indicative of the first attack strategy, wherein the evaluating includes:

for each narrative point of the first attack strategy, determining whether the received communication includes the presence of one or more terms matching to the group of terms corresponding to the narrative point;

determining whether the received communication includes the presence of terms matching to a threshold number of narrative points; and

determining whether terms present in the communication that match to the narrative points of the first attack strategy occur in a particular order; and

classifying the received communication based at least in part on the evaluation.

21. The system of claim **15** wherein the reviewer is associated with a reputation.

* * * * *

Attachment (2)

Proofpoint Email Protection

Detect and Block Both Malicious and Malware-less Email Threats

KEY BENEFITS

- Block BEC scams, phishing attacks and advanced malware at entry
- Raise user awareness with email warning tag
- Improve productivity with fast email tracing and email hygiene
- Scale for large enterprises with complete flexibility
- Provide operational efficiencies through automation of security operation and threat response
- Extend protection with integrated email authentication, email encryption, email DLP, Targeted Attack Protection and more
- Deliver industry-leading SLAs when deployed in the cloud:
 - 99.999% service availability
 - 100% virus protection
 - less than one-minute email latency
 - 99% blocked or redirected spam

Proofpoint Email Protection helps secure and control your inbound and outbound email. It uses machine learning and multilayered detection techniques to identify and block malicious email. It also dynamically classifies today's threats and common nuisances. And it gives you granular control over a wide range of email. This includes imposter email, phishing, malware, spam, bulk mail and more. It also offers complete flexibility with custom security policies and mail routing rules. It's also the most deployed email security solution by the Fortune 1000. And it scales for even the largest enterprise. What's more, it supports cloud, on-premises and hybrid installations.

Email is the No. 1 threat vector, with 96% of suspicious social actions arriving through email.¹ In addition to common email threats like phishing attacks and malware, emerging business email compromise (BEC) has posed a new threat to organizations. Email Protection catches both known and unknown threats that others miss. By processing billions of messages each day, Proofpoint sees more threats, detects them faster and better protects you against hard-to-detect malware-less threats, such as impostor email. With Email Protection, you can stop a vast majority of threats before they arrive in your user's inbox.

Catch Emerging Threats That Others Miss

Detect phishing, impostor and fraudulent email

Email Protection detects emerging threats before they can get to your user's inbox. Proofpoint Advanced BEC Defense powered by NexusAI is designed to effectively stop a wide variety of email fraud. That includes payment redirect and supplier invoicing fraud from compromised accounts. These types of threats require a more sophisticated detection technique, as there's often no malicious payload to detect.

Advanced BEC Defense is our ML and AI-powered detection engine. It is specifically designed to find and stop BEC attacks. It dynamically detects BEC by analyzing multiple

¹ Data Breach Investigations Report, Verizon, 2020.

message attributes. Some examples include:

- Message header data
- Sender's IP address (x-originating IP, reputation)
- Message body for urgency and words or phrases

It determines whether a message is a BEC threat. And it detects various BEC actor tactics. Such as:

- Reply-to pivots
- Use of malicious IPs
- Use of impersonated supplier domains

Advanced BEC Defense also provides granular visibility into BEC threat details. That includes BEC theme, gift card, payroll redirect, invoicing, lure or task. It provides observations about why the message was suspicious and message samples. That way, your security team can better understand and communicate about the attack. Data from NexusAI is then fed into the Proofpoint Nexus Threat Graph. It analyzes and correlates threat information across email, cloud, network and social from all of our customers. And thus gives you the protection to stay ahead of the threat landscape.

Block malicious and unwanted email

We've built multilayered detection techniques into Email Protection to defend against constantly evolving threats. With signature-based detection, it blocks known threats like viruses, trojan horses and ransomware. And it uses dynamic reputation analysis to continually assess local and global IP addresses to determine whether to accept email connections. Our unique email classifier also dynamically classifies a wide variety of emails. This includes impostor, phishing, malware, spam, bulk mail, adult content and circle of trust. And it quarantines incoming email by types. Together, these features help protect you at the first signs of malicious activity.

Track Down Any Email in Seconds

Email Protection has the most powerful search capability. With the smart search feature, you can easily pinpoint hard-to-find log data based on dozens of search criteria. You can also swiftly trace where emails come from and go to. Email Protection provides you with granular details of search results, including metadata with over a hundred attributes. The search is complete in seconds, not minutes.

You can download and export your search results by up to 1 million records. Moreover, several real-time reports are built into the product, giving you the detailed visibility into mail flow and trends. With this data, you can proactively address issues as they emerge.

Scales for Large Enterprise with Complete Flexibility

Email Protection supports the demands of the largest enterprises in the world. It allows you to create highly customizable email firewall rules at the global, group and user level. You can create any security policies and mail routing rules that fit your needs. And you can easily enforce them. Email Protection also provides the same benefits and greater flexibility with multiple deployment options. This includes on-premises hardware, virtual machine and SaaS.

Raise User Security Awareness

The email warning tag feature enables your users to make more informed decisions on the emails that fall into the gray area between clean and malicious. It surfaces a short description of the risk associated with a particular email. And it conveys the level of risk with different colors, which is easy to consume by your users. They can report suspicious email directly from the warning tag, even when they access email via mobile devices. This feature helps reduce the risk of potential compromise by making your users more cautious of uncertain email.

Email Protection also allows email admins to give users the ability to manage encrypted messages and low-priority emails like bulk mail, review quarantined messages and take actions directly in the Outlook task pane. User feedback is then transmitted to Proofpoint, helping us improve the global accuracy of bulk mail classification.

Centrally Manage across Email Encryption and DLP

You can easily extend your protection by adding Proofpoint Targeted Attack Protection, Email Fraud Defense, Email Encryption or Email Data Loss Prevention (DLP). While Email Protection provides you with basic email encryption and DLP capabilities, you can get more robust email encryption and DLP solutions through the same management console. This tight integration helps you manage sensitive data sent through email. It also prevents data leakage or data loss via email. And it satisfies several compliance requirements.

LEARN MORE

For more information, visit proofpoint.com.

ABOUT PROOFPOINT

Proofpoint, Inc. (NASDAQ: PFPT) is a leading cybersecurity and compliance company that protects organizations' greatest assets and biggest risks: their people. With an integrated suite of cloud-based solutions, Proofpoint helps companies around the world stop targeted threats, safeguard their data, and make their users more resilient against cyber attacks. Leading organizations of all sizes, including more than half of the Fortune 1000, rely on Proofpoint for people-centric security and compliance solutions that mitigate their most critical risks across email, the cloud, social media, and the web. More information is available at www.proofpoint.com.

©Proofpoint, Inc. Proofpoint is a trademark of Proofpoint, Inc. in the United States and other countries. All other trademarks contained herein are property of their respective owners. [Proofpoint.com](https://proofpoint.com)

Attachment (3)

Understanding BEC Scams: Supplier Invoicing Fraud

 proofpoint.com/us/blog/cybersecurity-essentials/understanding-bec-scams-supplier-invoicing-fraud

December 14, 2020

[Blog](#)

[Email and Cloud Threats](#)

Understanding BEC Scams: Supplier Invoicing Fraud

© 2020 Proofpoint, Inc. All rights reserved. Proofpoint and the Proofpoint logo are trademarks of Proofpoint, Inc. in the United States and other countries. Other trademarks are the property of their respective owners.



December 07, 2020 Tony Paterra

About the Series:

Business Email Compromise (BEC) and Email Account Compromise (EAC) afflict businesses of all sizes across every industry. More money is lost to this type of attack than any other cybercriminal activity. The FBI reported that from June 2016 to June 2019, companies reported \$26.2B in losses. And in 2019 alone, BEC scams accounted for more than half of all cybercrime losses—an estimated \$1.77B. The average loss per BEC incident in 2019 was \$74,723.

An indication of how pervasive a problem BEC/EAC is: Proofpoint blocks over 15,000 BEC/imposter messages a business day or nearly 4 million messages a year.

In this series, we cover BEC/EAC attack types that Proofpoint considers the most important for a business to be aware of. In each of these posts, we explain what the attack is and how it works so that you can better understand these BEC/EAC attacks and can better protect yourself. In this post, we focus on BEC/EAC supplier invoicing attacks.

What Is BEC Supplier Invoicing Fraud

BEC supplier invoicing scams are sophisticated and complex schemes to steal money by either presenting a fraudulent invoice as legitimate or by re-routing the payment to a bank account controlled by the attacker. When you consider the large dollars associated with supplier invoices, these scams are often the costliest for victim organizations. Proofpoint has stopped multiple supplier invoicing attempts where each incident was millions of dollars.

© 2020 Proofpoint, Inc. All rights reserved. Proofpoint and the Proofpoint logo are trademarks of Proofpoint, Inc. in the United States and other countries. Other trademarks are the property of their respective owners.

Proofpoint stopped multiple million dollar supplier invoicing scams in 2020.

BEC supplier invoicing fraud can be so successful that even prominent, well-known individuals can fall for them. In February 2020, Shark Tank's Barbara Corcoran nearly lost close to US\$400,000 to a BEC supplier invoicing attack. Happily for her, they were able to recover the funds before the attackers were able to collect them. That happy ending however is rare: BEC supplier invoicing scams rarely end with funds being successfully recovered.

This example is consistent with the high payoff successful supplier invoicing scams can yield. Proofpoint has successfully stopped invoicing fraud that could have yielded attackers millions of dollars if successful.

Similar to gift card scams and payroll diversion scams, supplier invoicing scams rely on social engineering and impersonation to convince the target victim to send money to the attackers. But what sets BEC supplier invoicing scams apart is not just the large dollar amounts often associated with these scams, but also the complex nature of these scams.

While gift card scams are relatively simple, using maybe one email targeting one employee, supplier invoicing scams are more byzantine involving compromise and impersonation of trusted vendors and carried out in multiple stages against multiple individuals and organizations. The impersonation can either be at an account level or at the domain level (e.g. domain lookalikes).

How BEC Supplier Fraud Works

Many of the BEC supplier invoicing attacks Proofpoint has observed indicate that these attacks originate from a legitimate email account that has been compromised. These compromised accounts are highly prized by threat actors. They can conduct extensive reconnaissance and fraudulent emails sent from the compromised account will pass email authentication controls (e.g. DKIM, SPF, DMARC) because they are sent from a legitimate account.

Once a legitimate transaction is identified, the threat actor “thread hijacks” an already in-progress email conversation about the transaction (step 3 in the diagram below). Since the attacker’s message is part of an email thread that the target victim reasonably believes to be legitimate, their message has greater credibility. As such requests for bank account changes due to audit or COVID-19 seem more plausible. This believability and trust are key elements of social engineering. By their very nature, thread hijacking attacks are very difficult, if not impossible for users to identify, making this a threat vector where technology countermeasures are particularly needed and useful.

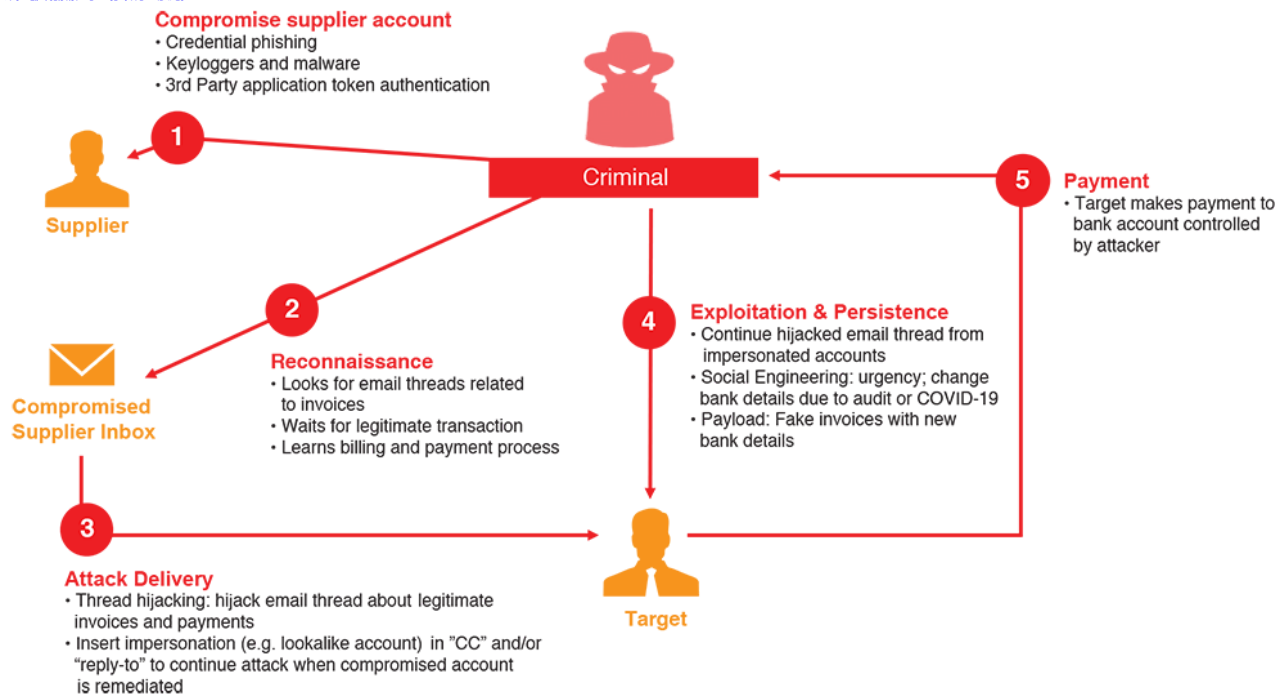


Figure 1: Anatomy of a Supplier Invoicing Fraud Attack

At this stage of the attack, the threat actor pivots to a supplier account impersonation tactic where the attacker inserts an impersonated account in the “reply-to” or “cc” of the email conversation. The impersonated accounts can be a lookalike of the supplier domain (e.g. supp1ier.com instead of supplier.com) or a lookalike of the account (e.g. janedoe[.]supplier[.]mail[.]com instead of jane[.]doe[.]supplier[.]com).

Why do this? The impersonation pivot allows the threat actor to maintain the email conversation with the target when the compromised account be remediated. In many cases, the email thread continues via the impersonated account. Shifting the conversation to the impersonated account also makes it more difficult for forensics and investigations because you lose the logs in the supplier SEG.

The images below show exchanges in a BEC supplier invoicing attack that Proofpoint stopped. The attackers in this particular attack started with the compromised supplier account where the attackers hijacked an existing email thread and utilized the

“impersonation reply-to pivot” tactic. They impersonated two separate trusted individuals and targeted two separate individuals in Finance and Accounting roles to try and make the fraud seem legitimate.

In Figure 2 we see an email message from the compromised account of a trusted individual from a supplier hijacking an existing email thread about several legitimate invoices. In this message the attackers are using the trusted individuals actual account and achieving impersonation through control granted by compromising that account.

CRM 1/31/21 04:00:33 CAC DOCUMENT 1/3 1/19/2021 1/19/2021 1/19/2021 1/19/2021 1/19/2021

From: Chris@supplier (compromised supplier account)
To: Jason (target)

External Message

Thanks Connie~

Dear Jason,

Hope you are well.

The following invoices are due or will be due in Apr. And now we haven't received the payment from you side. Could you please help to arrange the payment in Apr? Thank you.

Total amount: USD 2,791,867.92

Class	Number	Amount(USD)	Days Late	Transaction Date	Due Date
Invoice		179,976.49	43	12-26-2019	03-10-2020
Invoice		15,328.07	34	01-04-2020	03-19-2020
Invoice		36,128.50	31	01-07-2020	03-22-2020
Invoice		29,744.80	31	01-07-2020	03-22-2020
Invoice		62,243.65	29	01-09-2020	03-24-2020
Invoice		9,306.72	28	01-10-2020	03-25-2020
Invoice		8,846.00	28	01-10-2020	03-25-2020
Invoice		1,873.20	28	01-10-2020	03-25-2020
Invoice		3,439.44	27	01-11-2020	03-26-2020
Invoice		54,257.82	27	01-11-2020	03-26-2020
Invoice		1,267.58	24	01-14-2020	03-29-2020
Invoice		11,290.40	22	01-16-2020	03-31-2020

Figure 2: Hijacking Email Thread from Compromised Supplier Account to Target Account 1

In Figure 3 the attackers pivot from the compromised account to an impersonated account fully under the attacker's control. This account appears to be another trusted individual from the supplier sending a follow up email to both target victims. The impersonation is carried out through use of a fake account made to seem like the trusted individual's account.

From: Sobha_supplier@fake.com (impersonated supplier account)
To: Jason and Carrie (targets)

We have sent several emails to you as per below but still no response. Could you please kindly check below email from chris our finance and reply back. We need your urgent and positive respond to our below email. Be aware and noted that our company will not take or accept the blame if any payment is remitted to our old account and we cannot receive it. Hence to avoid problem in payments, We need you to reply to our emails. Please check below email from Chris and reply back immediately!!

(1) Regarding the payment for below invoices ([REDACTED]

Could you please confirm back how soon can the payment be settled to us? as we are urgently in need of funds we will appreciate if you can remit the payments to us before the end of this month. Please confirm back

(2) To avoid problem in payment, We will need to advise and provide to you with our Company offshore Bank account details for your safe remittance purpose so as not to have any error in receiving your payment. Please kindly reply back to us immediately so we can provide you the offshore Banking details for the remittance purpose

Figure 3: Email from Impersonated Account to Both Target Accounts

The fraud is made more credible by the fact that this second email references the same invoice numbers as the first email and both emails have been sent to one of the target victims. These two emails reinforce each other and that helps create a greater sense of legitimacy that the attackers hope will convince one of the target victims to respond and take action.

In Figure 4 the attackers have shifted back to the compromised supplier account and are requesting a change of bank to intercept the payment.

From Chris (compromised supplier account)
To: Carrie (target)

Dear Carrie

Thanks and noted your below email.

Hence as i have informed you in my below previous email, that we have been informed by our bankers that they are on Audit, there is no inflow or outflow of fund, Hence bank has put stop order to our company account until they solve the problem! Hence to protect our interest and the interest of our customers and to avoid any possible error in receiveig your payment., Our bankers has instructed that all payments from now onward is to be directed remitted into our company offshore bank account for safety receiving.

We will need to advise and provide to you with our Company offshore Bank account details for your safe remittance purpose so as not to have any error in receiving your payment. Please kindly reply back to us immediately so we can provide you the offshore Banking details for the remittance purpose

Thanks and awaiting your immediate response!! We will send you the bank details once we get your fast response

Figure 4: Email from Compromised Supplier Account requesting bank change for payment

It's also notable that this email is using both authority and urgency, both common social engineering tactics in BEC attacks. Also notable is that the fraudulent emails are devoid of any malware payload such as an attachment or URL. There are no links or attachments for the victims to click.

Taken as a whole, this shows how attackers weave together identity deception, authority, and urgency while using tactics like account compromise and impersonation pivot all to make a fraudulent bank account change request seem legitimate so that target will pay the invoices to the threat actor's bank account.

CRM 7/27/24-000333-CLC DOCUMENT 1-3 LPRR 08/04/27 LPRR 38 16 44 LPRR 16 17

How Proofpoint Protects Against BEC/EAC Supplier Invoicing Fraud Scams

Proofpoint provides a multi-layered solution to help organizations protect against supplier invoicing fraud. First, as part of Proofpoint Email Protection, NexusAI for BEC Detection dynamically analyzes a wide range of message attributes, including header information, domain, and message body to determine if a message is an impostor message. Proofpoint Email Protection delivers unique value to Proofpoint customers in the following ways:

- The Nexus Threat Graph aggregates and correlates trillions of threat data points across email, cloud accounts, domains and more. Threat visibility across multiple attack vectors is a critical element to identify a fraudulent message from a compromised account.
- Analyzing and blocking 15,000 BEC messages per day globally gives each Proofpoint customer an unfair advantage as NexusAI for BEC Detection has a large corpus of prior BEC attacks to compare message body contents against.
- Being in outbound mail flow enables us to understand bi-directional communications and identify whether there's an established history of communication or whether the message is from someone you've never communicated with but is impersonating a supplier.

Second, the Nexus Supplier Risk Explorer continuously maps your supply chain and uncovers what threats they may be sending to your organization. This visibility helps you understand which suppliers are the riskiest and allows you implement adaptive controls to mitigate that risk.

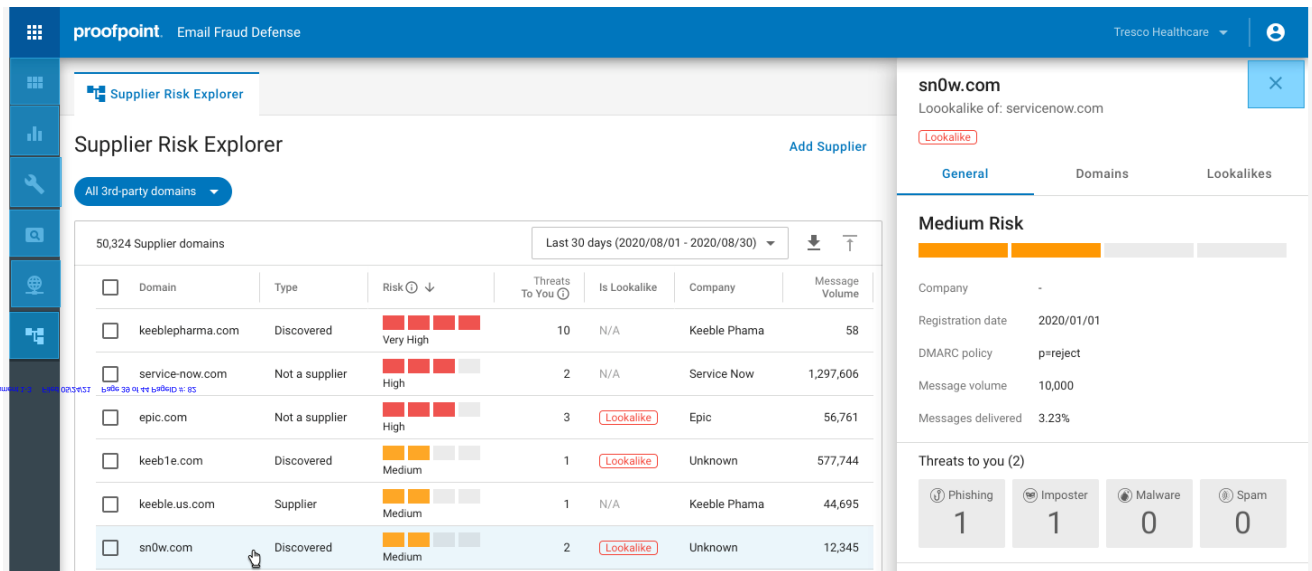
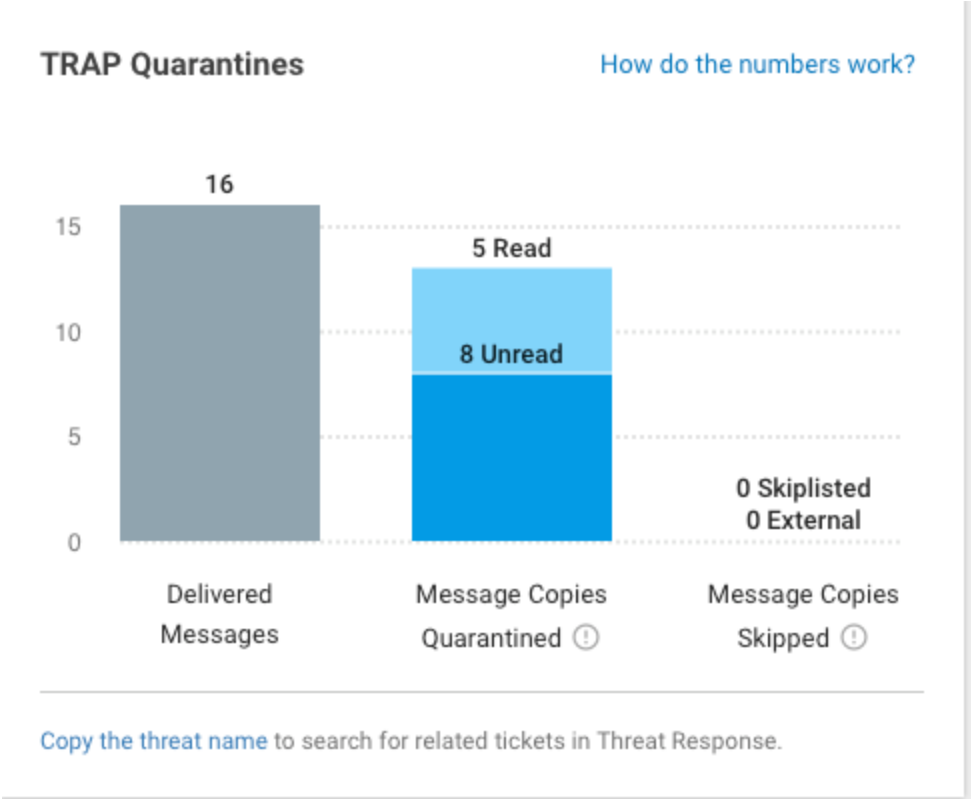


Figure 5: Nexus Supplier Risk Explorer automatically identifies your suppliers and the risks they pose

Third, Proofpoint provides incident responders with reporting to understand what type of BEC attack it is (e.g. invoicing, gift card, payroll) and who received it. This visibility is an important for prioritization efforts and informs what actions should be taken. Proofpoint also helps streamline incident response with the ability to automatically find and pull back delivered messages, included forwarded mail and distribution lists.




СРБГ Т-СТ-СА-00131-СБС ДОСЛУЖИЈ Т-3 ЕИЈБГ 0215415Т БРДБ 10 04 11 БРДБГД #: 83

СРБГ Т-СТ-СА-00131-СБС ДОСЛУЖИЈ Т-3 ЕИЈБГ 0215415Т БРДБ 10 04 11 БРДБГД #: 83

Attachment (4)

Using AI to Stop Threats and Reduce Compliance Risk

 proofpoint.com/us/blog/compliance/using-ai-stop-threats-and-reduce-compliance-risk

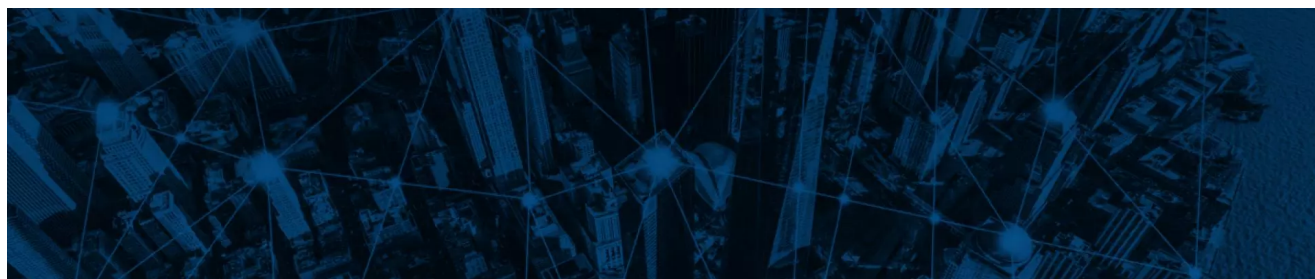
January 13, 2021

Blog

Archiving and Compliance

Using AI to Stop Threats and Reduce Compliance Risk

© 2021 Proofpoint, Inc. All rights reserved. | Privacy Policy | Terms of Service



January 13, 2021 Dan Rapp

Today's threat landscape is characterized by attackers preying on human vulnerability. Proofpoint research shows that nearly 99% of all threats require some sort of human interaction. Whether it is malware-free threats such as the different types of Business Email Compromise (BEC) or Email Account Compromise (EAC) like payroll diversion, account takeover, and executive impersonation, or malware-based threats like the Ryuk Ransomware that is carried by Bazalloader, people are falling victim to these attacks day-in and day-out.

Leveraging Machine Learning for a People-Centric Problem

To stop these types of attacks, organizations need to deploy a solution that can stay ahead of the ever-changing landscape and adapt to the way humans act. Machine Learning (ML) is a critical component in a robust cybersecurity detection strategy. It's faster and more effective than manual analysis and can quickly adapt to new and evolving threats and trends.

ML techniques aren't anything new to Proofpoint. Proofpoint NexusAI and ML techniques have been part of the Proofpoint DNA for over 10 years, and our current approach represents the third generation of our applied machine learning. We leverage a number of artificial intelligence cybersecurity techniques, which have become foundational in our efforts to help our customers better protect their people.

Applying NexusAI to Stop Threats

Proofpoint NexusAI employs many techniques to detect numerous targeted threats such as BEC, EAC, phishing, cloud attacks, multi-stage attacks, and others. As an example, BEC supplier invoicing fraud attacks are sophisticated and complex schemes to steal money by either presenting a fraudulent invoice as legitimate or by re-routing the payment to a bank account controlled by the attacker. It is extremely hard for traditional systems to detect due to two factors: they are very targeted and contain no payload. We are successful at stopping these threats with NexusAI for BEC Detection which uses ML to dynamically analyze a wide range of message attributes, including header information, domain, and message body to determine if a message is an impostor message. And the effectiveness of the ML models used in NexusAI for BEC detection is fueled by the dynamic, global corpus of BEC attack data we analyze on a daily basis – typically over 15,000 BEC messages.

Another example, credential phishing attacks, show inherent weakness as they try to mimic the look and feel of an organization's log-in page. This may be difficult to detect by analyzing manually, but with NexusAI, we leverage computer vision to detect and prevent emails pointing to that site.

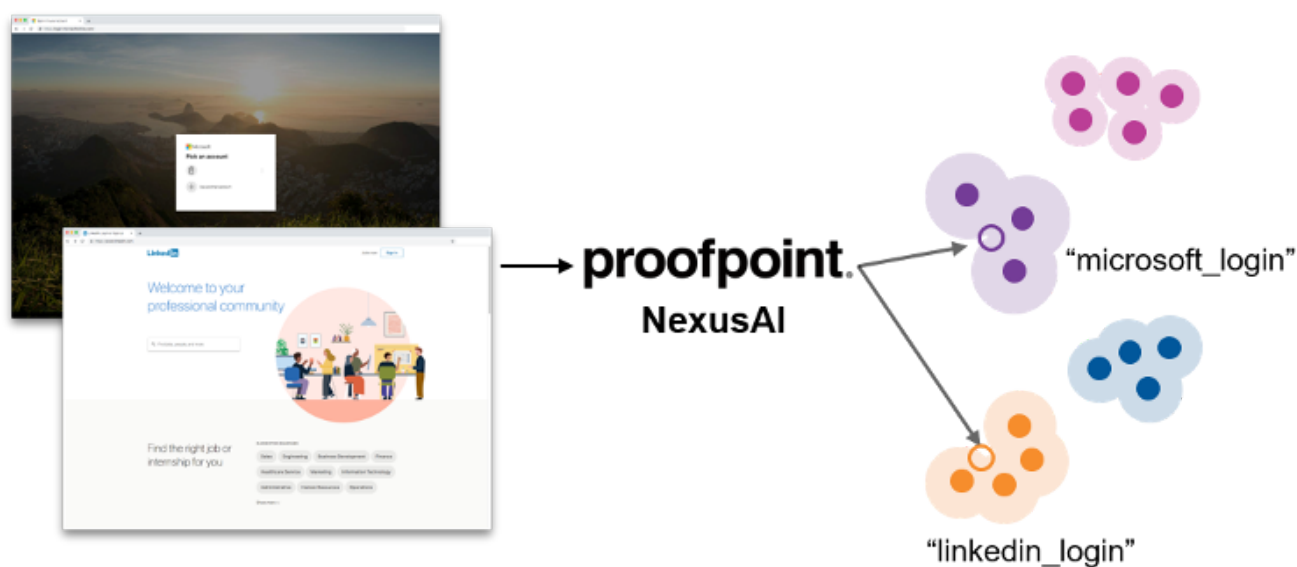


Figure 1. Proofpoint NexusAI – Computer Vision

We even leverage our NexusAI with our internal threat researchers. They use the NexusAI cybersecurity tools to better understand incidents as well as find like-incidents. This helps with improving the overall effectiveness of the entire Proofpoint security portfolio.

Data Dictates the Effectiveness of Machine Learning

ML systems are not like traditional software systems, behavior is derived from data and not hand-coded. That means ML systems are only as good as the data that trains them.

Proofpoint is unique in that we gather threat data from around the globe from leading enterprises in the F100, F1000, G2000, leading ISPs and hundreds of thousands of SMBs. We also gather data from across multiple attack vectors such as email, cloud, network and social media, which is critical as attackers augment their arsenal beyond email to include cloud applications.

- 2B+ daily emails scanned
- 22M+ cloud accounts monitored
- 6k+ worldwide IDS sensors
- 400k+ unique daily malware sample
- 400M+ domains monitored daily
- 86k+ social media accounts monitored
- 100+ threat actors tracked

This data is gathered from across industries and geographies is fed into the Nexus Threat Graph to learn everything we can about threats and train our NexusAI and detection technology. It's also used by our global team of threat researchers to analyze threats.

Everything learned from the data we see across our customer base is then fed back into Proofpoint products. This goes beyond only leveraging the data from a specific customer in our ML models but expands across our entire customer base. This data sharing helps us better identify campaigns and attack patterns we see across the threat landscape to provide the best cybersecurity and visibility for our customers to stay ahead of the threat landscape.

To learn more about how Proofpoint can help protect your people from threats [click here](#).

Subscribe to the Proofpoint Blog